

# Intermediate Machine Learning

## Lecture 1: Sparse Regression

Regression typically involves  $n$  pairs  $(x_1, y_1), \dots, (x_n, y_n)$

$y_i \in \mathbb{R}$  is the response

$x_i \in \mathbb{R}^p$  is the covariate or feature vector

The goal is to predict new  $Y$  from a given  $X$

Linear regression

$$m(x) = \mathbb{E}(Y|X)$$

Regression function

$$m(x) = \mathbb{E}(Y|X=x) = \int y P(Y|X=x) dx$$

Prediction Error

$$R = \mathbb{E}(Y - \hat{Y})^2 \text{ where } \hat{Y} \text{ is the prediction based on training data}$$

Bias-Variance Decomposition

$$R(\hat{m}) = \mathbb{E}(Y - \hat{m}(X))^2 = \int \text{bias}^2(x) p(x) dx + \int \text{Var}(x) p(x) dx + \sigma^2$$

$$\text{bias} = \mathbb{E}(\hat{m}(x)) - m(x)$$

$$\text{Var} = \text{Variance}(\hat{m}(x))$$

$$\sigma^2 = \mathbb{E}(Y - m(X))^2 \leftarrow \text{noise is intrinsic}$$

good prediction require a balance between bias and variance

Error can be decomposed into approximation error and estimation error

$$R(\hat{f}) - R^* = \underbrace{R(\hat{f}) - \inf_{f \in \mathcal{F}} R(f)}_{\substack{\text{estimation error} \\ \approx \text{variance}}} + \underbrace{\inf_{f \in \mathcal{F}} R(f) - R^*}_{\substack{\text{approximation} \\ \text{error} \\ \approx \text{bias}^2}}$$

Goal is to fit the model

$$m(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p = \beta^T x$$

Don't assume true regression function is linear

minimize training error

$$\frac{1}{n} \sum_{i=1}^n (y_i - \beta^T x_i)$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$X = \begin{matrix} & \begin{matrix} x_{11} & \dots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{matrix} \\ \begin{matrix} n \times p \end{matrix} & \end{matrix}$$

$$\hat{Y} = \hat{m}(X) = \hat{\beta}^T X$$

Improve the model to deal with high dimensions and more flexible predictions

Lecture 2: Sparse Regression continued

Ridge Regression

penalizes large coefficients

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (y_i - \beta^T x_i)^2 + \frac{\lambda}{n} \|\beta\|_2^2$$

$$\|\beta\|_2 = \sqrt{\sum_j \beta_j^2}$$

$$\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$$

$\lambda = 0 \Rightarrow$  least squares

$\lambda = \infty \Rightarrow \hat{\beta} = 0$  (variance approaches 0)

Use leave-one-out cross-validation

1. leave out  $(x_i, y_i)$
2. Find  $\hat{\beta}$
3. predict  $y_i = \hat{y}_{(-i)} = \hat{\beta}^T x_i$
4. Repeat for each  $i$

Computational shortcut for bias correction

$$\hat{R}(\lambda) \approx R_{\text{training}} + \frac{2P_{\lambda} \hat{\sigma}^2}{n}$$

$$H = X (X^T X + \lambda I)^{-1} X^T$$

$$P_{\lambda} = \text{trace}(H)$$

Ridge regression doesn't zero out any coefficients

We can use sparsity and convex relaxation to simplify the model selection process

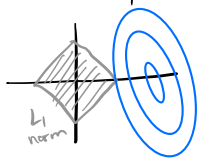
Sparsity: only a few predictors are relevant

Convex Relaxation: Replace model search with easier problem

$$\|\beta\|_q = (|\beta_1|^q + \dots + |\beta_p|^q)^{1/q}$$

L-2 norm doesn't measure sparsity but L-1 and L-0 do

Relates to convexity of the norms (L-1 and L-0 are convex)



overlap tends to occur at corners which zero out certain coefficients

Lasso Regression

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2n} \sum_{i=1}^n (y_i - \beta^T x_i)^2 + \lambda \|\beta\|_1$$

$$\|\beta\|_1 = \sum_j |\beta_j|$$

← No closed form

Penalized regression requires standardized coefficients

Risk estimation in the lasso

1. Find  $\hat{\beta}(\lambda)$  and  $\hat{S}(\lambda)$  for each  $\lambda$
2. Choose  $\lambda$  to minimize estimated risk
3. Let  $\hat{S}$  be the selected variables
4. Let  $\hat{\beta}$  be the least squares estimator using only  $\hat{S}$
5. Prediction  $\hat{y} = X^T \hat{\beta}$

# Lecture 3: Lasso, Smoothing and kernels

## Selecting $\lambda$

Re-fit the model with non-zero coefficients and apply leave-one-out cross validation

$$\hat{R}(\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_{(i)})^2 = \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \tilde{y}_i)^2}{(1 - H_{ii})^2} \approx \frac{1}{n} \frac{RSS}{(1 - \frac{1}{n})^2}$$

H is hat matrix  $S = \# \hat{\beta}_j \neq 0$

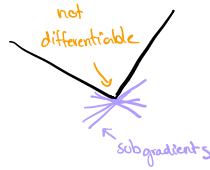
Choose  $\lambda$  to minimize risk

Algorithm for calculating Lasso

Consider  $L = \frac{1}{2}(Y - \beta)^2 + \lambda |\beta|$

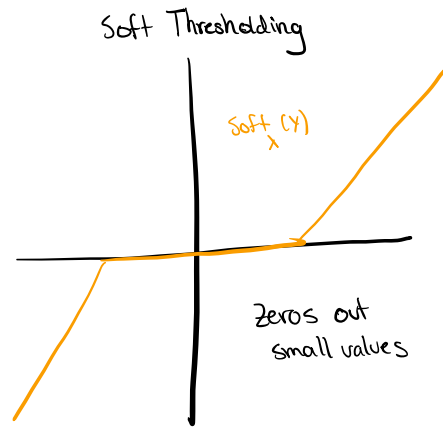
$$\frac{\partial L}{\partial \beta} = \beta - Y + \lambda v = 0$$

subgradient vector



$$\beta = Y - \lambda v = \begin{cases} Y - \lambda & \beta > 0 \\ Y + \lambda & \beta < 0 \\ Y - \lambda \frac{Y}{\lambda} & \beta = 0 \end{cases}$$

$$\Leftrightarrow \hat{\beta} = \begin{cases} Y - \lambda & Y > \lambda \\ Y + \lambda & Y < -\lambda \\ 0 & |Y| \leq \lambda \end{cases}$$



Compute least squares estimator and then run it through the soft thresholding function

## Lasso by coordinate descent

1. set  $\hat{\beta} = (0, \dots, 0)$  then iterate until convergence
2. For  $j = 1, \dots, P$ 
  - set  $r_i = y_i - \sum_{s \neq j} \hat{\beta}_s x_{si}$
  - set  $\hat{\beta}_j$  to be least squares fit of  $r_i$ s on  $x_j$
  - $\hat{\beta}_j \leftarrow \text{Soft}_{\lambda_j}(\hat{\beta}_j)$  where  $\lambda_j = \frac{\lambda}{\frac{1}{n} \sum x_{ij}^2}$
3. Use least squares  $\hat{\beta}$  on selected subset  $\hat{S}(\lambda)$

## Nonparametric Regression

Two types of kernels

1. Smoothing kernels: local averaging
2. Mercer kernels: Regularization

Smoothing kernel Estimator

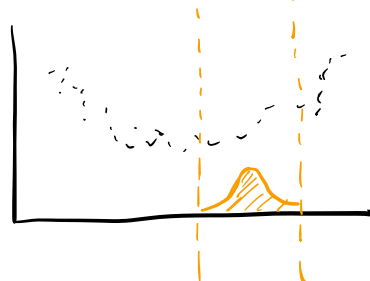
$$\hat{m}_h(x) = \frac{\sum_{i=1}^n y_i K_h(x_i, x)}{\sum_{i=1}^n K_h(x_i, x)}$$

← weights  
← normalization constant

$$K_h(x, z) = \exp\left(-\frac{\|x - z\|^2}{2h^2}\right)$$

← gaussian kernel  
h is the bandwidth

Weight points according to the kernel



Small  $h \rightarrow$  large variance, large  $h \rightarrow$  large bias

$$\text{Risk} = \mathbb{E}(Y - \hat{m}_h(X))^2 = \text{bias}^2 + \text{variance} + \sigma^2$$

under mild assumptions

$$\text{bias}^2 \approx h^4$$

$$\text{variance} \approx \frac{1}{nh^p}$$

The curse of dimensionality

risk goes down no faster than  $n^{-4/(4+p)}$

To make risk of size  $\epsilon$  we need

$$n \geq \left(\frac{1}{\epsilon}\right)^{4+p/4} \leftarrow \text{sample size grows exponentially}$$

We choose  $h$  by estimating risk  $R(h)$

1. Omit  $(x_i, y_i)$  to get  $\hat{m}_{h(-i)}$  then predict  $\hat{y}_{(-i)} = \hat{m}_{h(-i)}(x_i)$

2. Repeat for all observations

3. Cross validation risk estimate

$$\begin{aligned} \hat{R}(h) &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_{(-i)})^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{y}_i}{1 - L_{ii}} \right)^2 \\ L_{ii} &= \frac{K_h(x_i, x_i)}{\sum_{i=1}^n K_h(x_i, x_i)} \end{aligned}$$

## Additive Models

compromise between linear models and nonparametric models

$$\hat{m}(X) = \hat{m}_1(x_1) + \dots + \hat{m}_p(x_p)$$

Each  $\hat{m}_j(x_j)$  is estimated by smoothing while holding other functions fixed



# Lecture 4: Smoothing and Density Estimation

## Smoothing Kernel

$$\hat{m}_h(x) = \frac{\sum_{i=1}^n Y_i K_h(X_i, x)}{\sum_{i=1}^n K_h(X_i, x)}$$

← local weighted average of  $Y_i$  near  $x$

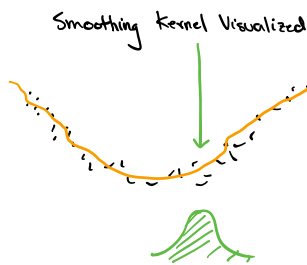
$$K_h(x, z) = \exp\left(-\frac{\|x-z\|^2}{2h^2}\right)$$

$h > 0$  is the bandwidth

$h$  controls bias variance tradeoff

Small  $h \rightarrow$  high variance

Large  $h \rightarrow$  high bias



## Estimating Risk with Cross-Validation

① Omit  $(X_i, Y_i)$  to get  $\hat{m}_h(\cdot)$  and predict  $\hat{Y}_{(i)} = \hat{m}_h(\cdot)(X_i)$

② Repeat for all observations

③ Cross-validated estimate of risk

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_{(i)})^2$$

Shortcut formula

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \frac{(Y_i - \hat{Y}_{(i)})^2}{1 - L_{ii}}$$

$$L_{ii} = \frac{K_h(X_i, X_i)}{\sum_{t=1}^n K_h(X_i, X_t)}$$

Proof  $\rightarrow$

$$\hat{Y}_{(i)} = \frac{\sum_{j \neq i} K_{ij} Y_j}{\sum_{j \neq i} K_{ij}} = \frac{\sum_j K_{ij} Y_j - K_{ii} Y_i}{\sum_j K_{ij} - K_{ii}} = \frac{\sum_j L_{ij} Y_j - L_{ii} Y_i}{1 - L_{ii}} = \frac{Y_i - L_{ii} Y_i}{1 - L_{ii}} \Rightarrow (Y_i - \hat{Y}_{(i)})^2 = \left(Y_i - \frac{Y_i - L_{ii} Y_i}{1 - L_{ii}}\right)^2 = \left(\frac{Y_i - \hat{Y}_{(i)}}{1 - L_{ii}}\right)^2$$

where  $K_{ij} = K_h(X_i, X_j)$

## General procedure for choosing bandwidth

- ① Compute  $\hat{m}_h$  for each  $h$
- ② Estimate risk  $\hat{R}(h)$  using LOOCV
- ③ Choosing  $\hat{h}$  to minimize  $\hat{R}(h)$
- ④ Let  $\hat{m}(x) = \hat{m}_{\hat{h}}(x)$

## Kernel Density Estimation

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(X_i, x) = \frac{1}{n} \cdot \frac{1}{h^p} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right)$$

↑ why

We require  $\int K(u) du = 1$  and  $K \geq 0$  is symmetric around 0

## Bias-variance tradeoff

$$\text{bias}^2(x) \approx h^4$$

$$\text{var}(x) \approx \frac{1}{nh^p}$$

## Derivation

$$\hat{f}_h(x) = \frac{1}{nh^p} \sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)$$

$$\text{Bias} : \mathbb{E} \hat{f}_h(x) - f(x)$$

$$\mathbb{E} \hat{f}_h(x) = \frac{1}{nh^p} \sum_{i=1}^n \mathbb{E} \left( K\left(\frac{X - X_i}{h}\right) \right) = \frac{1}{h^p} \mathbb{E} \left( K\left(\frac{X_i - x}{h}\right) \right) = \frac{1}{h^p} \int K\left(\frac{u-x}{h}\right) f(u) du = \int K(r) f(x+hr) dr$$

true density  
↓

change of variable

$$\frac{u-x}{h} = r \quad u = x+hr \quad du = h^p dr$$

$$= \int K(r) \left\{ f(x) + hr^T \nabla f(x) + \frac{h^2}{2} r^T \nabla^2 f(x) r + O(h^3) \right\} dr$$

→ Taylor expansion around  $x$  as  $h \rightarrow 0$

$$= \underbrace{f(x) \int K(r) dr}_1 + \underbrace{\nabla f(x)^T h \int r K(r) dr}_0 + \frac{h^2}{2} \int K(r) r^T \nabla^2 f(x) r dr + O(h^3) \leq ch^2$$

expected value of even function

$$\therefore \text{bias} \leq ch^2$$

Variance:  $\mathbb{E}[\hat{f}(x)^2] - \mathbb{E}[\hat{f}(x)]^2$

$$\mathbb{E}[\hat{f}(x)^2] = C_2 \frac{n}{n^2 h^{2p}} \int K\left(\frac{u-x}{h}\right)^2 f(u) du = C_2 \cdot \frac{1}{nh^p} \int K(r)^2 f(x+hr) dr = \frac{C_2}{nh^p} f(x) \int K(r)^2 dr + O\left(\frac{1}{nh^p}\right)$$

→ Same change of variables
→ Taylor approx

$\therefore \text{var} \leq \frac{1}{nh^p}$

Risk  $\leq C_1 h^4 + C_2 \frac{1}{nh^p}$

Heuristically, we can set  $h^4 = \frac{1}{nh^p}$

Risk  $\approx \left(\frac{1}{n}\right)^{4/4+p}$

$h = \frac{1}{n^{4/4+p}}$

Sample size grows exponentially as we increase dimensions to maintain the same risk levels

Kernel smoothing estimator is a plug-in estimator of the kernel density estimator subject to the same curse of dimensionality

**Lecture 5: Mercer kernels**

At a high level Mercer kernels are penalties to smoothing

A kernel is a bivariate function  $K(x, x')$  that can be used as a measure of similarity between  $x$  and  $x'$

A Mercer kernel has a special property that for  $n$  points the  $n \times n$  matrix where

$K_{ij} = [K(x_i, x_j)]$  is positive semi-definite ← necessary and sufficient condition

↑  
no negative eigenvalues

Gaussian and triangular kernels are Mercer kernels

Mercer kernel  $K(x, x')$  is symmetric and positive semi-definite bivariate function

$$\iint f(x)f(x') K(x, x') dx dx' \geq 0$$

for all univariate functions  $f$

We can create a set of basis functions based on  $K$  by fixing  $z$

$k(z, x) = K_z(x)$

Defining a norm from the kernel

Suppose  $f(x) = \sum_r \alpha_r K_z(x)$  and  $g(x) = \sum_s \beta_s K_{z_0}(x)$

$\langle f, g \rangle_K = \sum_r \sum_s \alpha_r \beta_s K(z_r, z_0)$

inner product →  $= \alpha^T K \beta$

$K = [K(z_r, z_0)]$

Positive semidefinite condition ensures that the inner product is positive

$$\|f\|_K^2 = \langle f, f \rangle_K = \sum_r \sum_s \alpha_r \alpha_s K(z_r, z_0) = \alpha^T K \alpha \geq 0$$

## Reproducing Kernel Hilbert Space

$$\langle f, k_x(\cdot) \rangle_K = f(x)$$

reproduces the function through the inner product

Norm allows us to penalize overly complex functions

$$\text{Minimize} \sum_{i=1}^n (y_i - m(x_i))^2 + \lambda \|m\|_K^2 \leftarrow \text{ridge regression}$$

over RKHS

## Representer Theorem

Let  $\hat{m}$  minimize

$$J(m) = \sum_{i=1}^n (y_i - m(x_i))^2 + \lambda \|m\|_K^2$$

Then

$$\hat{m}(x) = \sum_{i=1}^n \alpha_i K(x_i, x)$$

for some  $\alpha_1, \dots, \alpha_n$

Plugging  $\hat{m}$  into  $J$  we get

$$J(\alpha) = \|Y - K\alpha\|^2 + \lambda \alpha^T K \alpha \leftarrow \text{kernel ridge regression}$$

to minimize  $J$

$$\hat{\alpha} = (K + \lambda I)^{-1} Y$$

$$\hat{m}(x) = \sum_i \hat{\alpha}_i K(x_i, x)$$

Estimator depends on regularization parameter  $\lambda$

We choose via cross validation

## Smoothing kernels vs. Mercer kernels

smoothing kernels: bandwidth  $h$  controls the amount of smoothing

Mercer kernels: norm  $\|f\|_K$  controls the amount of smoothing

Both subject to the curse of dimensionality

## Kernels from features

If  $x \rightarrow \varphi(x) \in \mathbb{R}^d$ , we can define a Mercer kernel

$$K(x, x') = \varphi(x)^T \varphi(x')$$

Conversely we can use the spectral theorem to get a feature map from a Mercer kernel

## Lecture 6: Mercer kernels and Neural Networks

### Basics of Neural Nets

In linear regression our loss function is given by

$$\begin{aligned} \mathcal{J}(x, y) &= \frac{1}{2} (y - B^T x - B_0)^2 \\ &= \frac{1}{2} (y - f(x))^2 \end{aligned}$$

$$f(x) = B^T x + B_0$$

If we add a layer

$$\mathcal{J} = \frac{1}{2} (y - f(x))^2$$

$$f(x) = B^T h(x) + B_0 \text{ where } h(x) = Wx + b$$

However this model is still linear, just reparameterized

Adding non-linearities

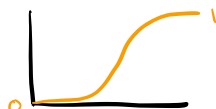
$$h(x) = \varphi(Wx + b) \leftarrow \text{Activation threshold}$$

Common Non-linearities

$$\tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$

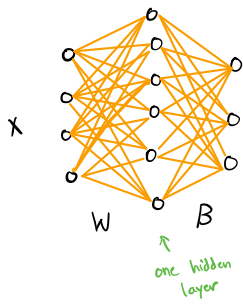
$$\text{sigmoid}(u) = \frac{e^u}{1 + e^u}$$

$$\text{relu}(u) = \max(u, 0)$$



$\leftarrow$  piecewise linear boundaries

## Multi-layer Perception



Neural Networks are just parametric regression models with a restricted type of non-linearity

## Backpropagation

Calculate derivatives via chain rule starting from the final layer

## Classification

To calculate final probabilities we use

$$(P_1, P_2, P_3) = \frac{1}{e^{f_1} + e^{f_2} + e^{f_3}} (e^{f_1}, e^{f_2}, e^{f_3})$$

## Lecture 7: Neural Networks and Convolutional Neural Networks

Neural Nets without non-linear components are simply least squares or logistic regression

Stochastic Gradient Descent

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} \text{Loss}(\Theta)$$

← step size

Compute gradients via chain rule

work your way backwards from the final node

Neural Nets are robust to overfitting

## Double Descent

$$f(x) = B^T \psi(Wx)$$

← inputs  
← first layer weights  
activation function

$$x \in \mathbb{R}^d$$

$$W \in \mathbb{R}^{p \times d}$$

If we fix  $W$  to be gaussian weights

$$\psi(Wx) = h(x)$$

← random features

Linear model in the features

$(x_i, y_i)$  for  $i=1, \dots, n$

$$X_i = h(x_i) \in \mathbb{R}^p$$

Consider  $(X_i, y_i)$

$$X \in \mathbb{R}^{n \times p}$$

$$X_i = X^T$$

Classical Regime ( $n > p$ )

$$\text{OLS: } \hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\hat{y} = X \hat{\beta}$$

Overparameterized Regime ( $n < p$ )

OLS doesn't exist

Minimum norm solution

$$\hat{\beta}_{\text{mn}} = X^T (X X^T)^{-1} Y \in \mathbb{R}^p$$

$$\hat{y} = X \hat{\beta}_{\text{mn}} = X X^T (X X^T)^{-1} Y = Y$$

Functions to interpolate

$$\mathbb{E} \hat{y}_i = \text{max}_i?$$

$$(\beta - \hat{\beta}_{min})^T \hat{\beta}_{min} = (\beta - \hat{\beta}_{min})^T X^T (XX^T)^{-1} Y$$

$$= (X(\beta - \hat{\beta}_{min}))^T (XX^T)^{-1} Y = 0$$

$$\|B\|^2 = \|\beta - \hat{\beta}_{min} + \hat{\beta}_{min}\|^2$$

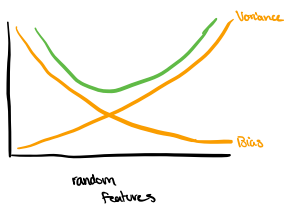
$$= \|\hat{\beta}_{min}\|^2 + \|\beta - \hat{\beta}_{min}\|^2 \geq \|\hat{\beta}_{min}\|^2$$

Ridge Regression

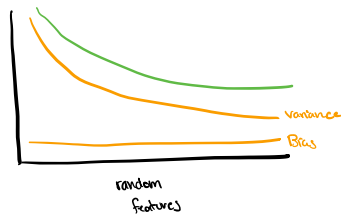
$$\hat{\beta}_\lambda = (X^T X + \lambda I_p)^{-1} X^T Y \xrightarrow{\lambda \rightarrow 0} \begin{cases} \text{OLS} & \text{Classical regime} \\ \text{minimum norm} & \text{Overparameterized regime} \end{cases}$$

Bias-Variance Tradeoff

Classical Regime



Overparameterized Regime



Kernel Connection

Consider the neural net as a function

$$f(\theta(x))$$

Neural Tangent Kernel  $\rightarrow$  Associated kernel for neural net

$$K(x, x') = \nabla_{\theta} f(x)^T \nabla_{\theta} f(x')$$

derivatives of the function

In the random features model

$$\nabla_{\theta} f_{\theta}(x) = \nabla_{\beta} R^T h(x) = h(x) \in \mathbb{R}^p$$

$$K = HH^T \leftarrow \text{Mercer Kernel}$$

## Lecture 8: Convolutional Neural Networks

Convolution sweeps a function  $K(x)$  over the values of another function

$$K * f(x) = \int K(x-u) f(u) du$$

Same operation as smoothing kernels

Convolution can be considered a linear operation

Filters can be used as an edge detector

Principles of CNN

- Parameter Tying: let neurons share the same weights
- Kernel Learning: kernels should not be fixed, instead they should be learned
- Invariance: Doesn't matter where features are selected

CNN is most commonly used with images or spectral data

Tensor is an n-dimensional array

Convolutional Layer

- Set of kernels are swept over the input
- Kernels are tensors with trained weights
- Resulting convolution is the feature map
- Equivalent to a traditional layer but weights are constrained to be sparse and tied

Adding nonlinearity

hit the featuremap with an intercept and an activation function

Convolutional layers are intended to learn "features" of the input

Pooling the feature map reduces dimension

Max pooling: max value in a block

Average pooling: Average value in a block

Compounded convolutional layers can learn features of features

First layer tends to be edge detection

Second layer builds shapes from edges

Backpropagation

Forward pass builds predictions

Backward pass propagates predictive error to update weights

Typically each kernel has the same number of channels as the input tensor

Each kernel outputs a feature map

Recombine feature maps to channels

# channels in output is # kernels in the layer

First dimension of input output tensor is always n

Dropout

Randomly zero out entries in a tensor

prevents overfitting

Flattening

Add dense layers after one or more convolutional layers

standard layers used to make the final prediction

## Lecture 9: Nonparametric Bayes: Gaussian Processes

Nonparametric Bayes

Extension of bayes inference to infinite dimensions

Typically neither prior nor posterior have a defined density

Well defined posterior

Statistical Problem	Frequentist Approach	Bayesian Approach
Regression	Kernel Smoother	Gaussian Process
CDF Estimation	Empirical CDF	Dirichlet Process
Density Estimation	Kernel Density estimation	Dirichlet Process mixture

Stochastic process is a collection of random variables with a shared text

Gaussian Process

$$Y_i = m(x_i) + \epsilon_i \quad i=1, \dots, n$$

Frequentist Kernel Estimator

$$\hat{m}(x) = \frac{\sum_{i=1}^n x_i K\left(\frac{x-x_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)}$$

↑ kernel      ↑ width

Bayesian method requires a gaussian prior

Exploits gaussian conditional/marginals

Slices of multivariate gaussians are gaussian with constant variance

A stochastic process is a gaussian process if for each set of points  $x_1, \dots, x_n$  the vector  $m(x_1), \dots, m(x_n)$  is normally distributed

$$(m(x_1), \dots, m(x_n))^T \sim \mathcal{N}(\mu(x), k(x))$$

↑ Mercer kernel

Gaussian Process Prior

$$\pi(m) = (2\pi)^{-n/2} |K|^{-1/2} \exp\left(-\frac{1}{2} m^T K^{-1} m\right)$$

Change of variable

$$\pi(\alpha) = (2\pi)^{-n/2} |K|^{-1/2} \exp\left(-\frac{1}{2} \alpha^T K \alpha\right)$$

Gaussian prior favors functions with weight on eigenvectors with large eigenvalues

### Smother functions

## Lecture 10: Nonparametric Bayes: Gaussian Processes + Dirichlet Processes

Maximum a posteriori estimate for gaussian process is equivalent to mercer kernel regression

Prior prefers smooth functions but the definition of smoothness is dictated by the mercer kernel (bandwidth)

If  $X_1$  and  $X_2$  are jointly gaussian, then the conditional distributions are also gaussian

Covariance of both distributions is independent of  $X$

Predicting a New Point

$$\text{Let } K = (K(X_1, X_{n+1}), \dots, K(X_n, X_{n+1}))$$

Then  $Y_1, \dots, Y_{n+1}$  are jointly gaussian with covariance

$$\begin{pmatrix} K + \sigma^2 I & K \\ K^T & K(X_{n+1}, X_{n+1}) + \sigma^2 \end{pmatrix}$$

The posterior mean and variance are as follows

$$\mathbb{E}(Y_{n+1} | X_{1:n}, Y_{1:n}) = K^T (K + \sigma^2 I)^{-1} Y$$

$$\text{Var}(Y_{n+1} | X_{1:n}, Y_{1:n}) = K(X_{n+1}, X_{n+1}) + \sigma^2 - K^T (K + \sigma^2 I)^{-1} K$$

↑  
select via  
domain knowledge

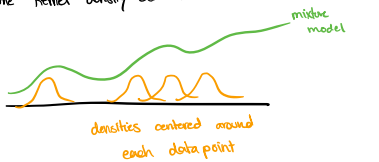
The Dirichlet Process is analogous to the gaussian process

method for density estimation instead of regression

Bayesian approach to density estimation

Every partition of sample space has a dirichlet distribution

Recall the kernel density estimate



### Building the Bayesian Model

Construct prior via imaginary data ← guess a beta distribution

Each sample from a Dirichlet prior has a random collection of weights assigned to a random selection of data

Each sample from a Dirichlet process mixture has a random collection of weights assigned to a random selection of model parameters

Stick Breaking Process for Dirichlet Process mixture

1. Draw  $\theta_1, \theta_2, \dots$  independently from  $F_0$

2. Draw  $V_1, V_2, \dots \sim \text{Beta}(1, \alpha)$

$$\text{Set } w_j = V_j \prod_{i=1}^{j-1} (1 - V_i)$$

3. Fit the model

$$f(x) = \sum_{j=1}^{\infty} w_j f(x | \theta_j)$$

A dirichlet process is a stochastic process where every partition of the sample space has a dirichlet distribution

We use an algorithm to sample from the posterior

No closed form

## Lecture 11: Approximate Inference: Simulation and Variational Methods

Simulation and Variational methods are ways to work with Bayesian posteriors

### Generative Model

1. Choose  $Z$  ← training point

2. Given  $Z$ , generate  $X$  ← gaussian sample

### Inverting generative Model

1. Given  $X$

2. What  $Z$  generated it?

### Ising Model

$Z_1, \dots, Z_n$  binary variables

$$P_{\beta}(z_1, \dots, z_n) \propto \exp\left(\sum_{s \in V} B_s z_s + \sum_{(s,t) \in E} J_{st} z_s z_t\right)$$

## Gibbs Sampler

1. Choose vertex  $s \in V$  at random

2. Sample  $u \sim \text{Uniform}(0,1)$

$$z_s = \begin{cases} 1 & u \leq \left(1 + \exp(-B_s - \sum_{t \in N(s)} B_{st} z_t)\right)^{-1} \\ 0 & \text{otherwise} \end{cases}$$

← akin to logistic

↑  
conditional probability  
 $z_s \in \{0,1\}$

3. Iterate

Encourages neighboring pixels to be similar

Gibbs sampling is a stochastic approximation

## Deterministic Approximation

### Mean Field Approximation

1. Choose vertex  $s \in V$  at random

2. Update  $\mu_s = \left(1 + \exp(-B_s - \sum_{t \in N(s)} B_{st} \mu_t)\right)^{-1}$

3. Iterate

## Mixture Model

Fix two distributions  $F_0$  and  $F_1$  with densities  $f_0(x)$  and  $f_1(x)$

$\theta \sim \text{Beta}(\alpha, \beta)$  ← mixing coefficient

$X|\theta \sim \theta F_1 + (1-\theta)F_0$

## Likelihood

$$p(x_{1:n}) = \int_0^1 \text{Beta}(\theta | \alpha, \beta) \prod_{i=1}^n (\theta f_1(x_i) + (1-\theta) f_0(x_i)) d\theta$$

## Gibbs Sampler

1. Sample  $z_i | \theta, x_{1:n}$

2. Sample  $\theta | z_{1:n}, x_{1:n}$

For step 1 sample  $u_i \sim \text{Uniform}(0,1)$

$$z_i = \begin{cases} 1 & u_i \leq \frac{\theta f_1(x_i)}{\theta f_1(x_i) + (1-\theta) f_0(x_i)} \\ 0 & \text{otherwise} \end{cases}$$

For step 2

$$\theta \sim \text{Beta}\left(\sum_{i=1}^n z_i + \alpha, n - \sum_{i=1}^n z_i + \beta\right)$$

## Lecture 12: Variational Inference and VAEs

### Inventing Models

Bayesian set up:

1. Choose  $\theta$
2. Given  $\theta$ , generate  $X$  sample

Posterior Inference:

1. Given  $X$
2. What  $\theta$  generated  $X$ ?

Approximate these answers via stochastic or deterministic methods

Gibbs Sampling → Stochastic

Variational Methods → Deterministic



## Stochastic Approximation (Gibbs Sampler)

Iterate until convergence

1. Choose vertex  $s \in V$  at random
2. Sample  $z_s$  while holding others fixed

$$\theta_s = \text{sigmoid} \left( B_s + \sum_{t \in N(s)} B_{st} z_t \right)$$

$$z_s | \theta_s \sim \text{Bernoulli}(\theta_s)$$

In Bayesian notation

1. Select a component  $s$
2. Calculate conditional distribution  $P(\theta_s | \theta_{\neq s}, X)$
3. Sample  $\rightarrow$

Repeat

Low dimensional so it is tractable calculation

## Deterministic Approximation (Mean Field variational)

Iterate until convergence

1. Choose vertex  $s \in V$  at random
2. Update mean  $M_s$  holding others fixed

$$M_s = \text{sigmoid} \left( B_s + \sum_{t \in N(s)} B_{st} M_t \right)$$

$$\frac{1}{1-e^{-x}}$$

## Comparison of Deterministic and Stochastic Approximation

- All  $z_s$  are random
- The  $M_s$  are deterministic
- Gibbs converges in distribution but mean field is numerical
- Gibbs approximates full distribution but mean field only approximate mean of each node

## A Finite Mixture Model

Consider two distributions  $F_0$  and  $F_1$  with densities  $f_0(x)$  and  $f_1(x)$  with the following mixture model

mixing parameter  $\rightarrow \theta \sim \text{Beta}(\alpha, \beta)$   
 $X | \theta \sim \theta F_1 + (1-\theta) F_0$

Likelihood for given data

$$P(X_{1:n}) = \int_0^1 \text{Beta}(\theta | \alpha, \beta) \prod_{i=1}^n (\theta f_1(x_i) + (1-\theta) f_0(x_i)) d\theta$$

## Gibbs Sampling Approach

1. Sample  $z_i | \theta, X_{1:n}$  for  $i=1, \dots, n$

repeat  $\left\{ \begin{array}{l} 1 \\ 0 \end{array} \right. \propto \theta f_1(x_i) \leftarrow P(z_i=1) = \frac{\theta f_1(x_i)}{\theta f_1(x_i) + (1-\theta) f_0(x_i)}$   
 $\propto (1-\theta) f_0(x_i)$

2. Sample  $\theta | z_{1:n}, X_{1:n}$

$$\theta \sim \text{Beta} \left( \sum_{i=1}^n z_i + \alpha, n - \sum_{i=1}^n z_i + \beta \right)$$

Sampled from a mixture of beta distributions over multiple iterations

## Collapsed Gibbs Sampler

Integrate out  $\theta$  and sample  $z_1 \dots z_n$

better algo

## Variational Inference (Deterministic Approximation)

Goal: compute  $P(\theta, z | X)$

Since the goal is too complicated, we approximate via  $q(\theta, z)$  that has a nice form  $\leftarrow$  fixed  $x$

$q$  is a function of variational parameters optimized for each  $x$

Maximize lower bound on  $P(X)$

evidence

Bound  $p(x)$  or alternatively  $\log p(x)$

$$\begin{aligned} \log p(x) &= \log p(x) \int q(\theta, z) d\mu(\theta, z) \\ &= \int q(\theta, z) \log p(x) d\mu(\theta, z) = \int q(\theta, z) \log \left( \frac{P(\theta, x, z)}{P(x, \theta, z)} \cdot \frac{q(\theta, z)}{q(\theta, z)} \right) d\mu(\theta, z) \end{aligned}$$

Rearrange

$$\begin{aligned} \log p(x) &= \int q(\theta, z) \log \left( \frac{P(\theta, x, z)}{P(\theta, z|x)} \cdot \frac{q(\theta, z)}{q(\theta, z)} \right) d\mu(x, z) \\ &= \underbrace{\int q(\theta, z) \log \left( \frac{q(\theta, z)}{P(\theta, z|x)} \right) d\mu(\theta, z)}_{\substack{\text{drop non-negative term} \\ \text{KL divergence}}} + \int q(\theta, z) \log \left( \frac{P(\theta, x, z)}{q(\theta, z)} \right) d\mu(\theta, z) \end{aligned}$$

$$\log p(x) \geq \int q(\theta, z) \log \frac{P(\theta, x, z)}{P(\theta, z)} d\mu(\theta, z)$$

$$= \mathbb{E}_q \log P(\theta, x, z) + H(q)$$

$\uparrow$  cross entropy  
 favors fit to the data  
 $\uparrow$  entropy  
 favors spreading out variational distribution

ELBO: Evidence Lower Bound

Numerically optimize w.r.t terms (variational parameters)

### Variational Algorithm for Mixture

Iterate for  $q_{1:n}$  and  $\gamma_1, \gamma_2$

1. Holding  $q_i$  fixed set

$$\gamma_1 = \alpha + \sum_{i=1}^n q_i \quad \gamma_2 = B + n - \sum_{i=1}^n q_i$$

2. Holding  $\gamma_1$  and  $\gamma_2$  fixed set

$$q_i = \frac{f_1(x_i) \exp \psi(\gamma_1)}{f_1(x_i) \exp \psi(\gamma_1) + f_2(x_i) \exp \psi(\gamma_2)}$$

After convergence

$$\hat{p}(\theta | x_{1:n}) = \text{Beta}(\theta | \gamma_1, \gamma_2)$$

Easier to use than gibbs sampling

Posterior is estimated as a single Beta distribution

## Lecture 12: Variational Autoencoders

Variational Inference: strategy

goal is to compute  $p(\theta, z | x)$  but that's too complicated

Instead, we approximate with  $q(\theta, z)$  that has a nice form  $\leftarrow$  typically gaussian

$q$  is a function of variational parameters that is optimized for each  $x$

Maximize a lower bound on  $p(x)$

ELBO is the lower bound on  $\log p(x)$

$$\log p(x) = H(q) + \mathbb{E}_q(\log p(x, z, \theta))$$

$\uparrow$  entropy

maximize over the parameters of  $q$

Entropy is maximized when the variational distribution is spread out

$$ELBO = H(q) + \mathbb{E}_q \log p(x, \theta, z)$$

$$H(q) + \mathbb{E}_q \log p(x|z, \theta) - \mathbb{E}_z \log p(z, \theta) \quad \leftarrow \text{Prior}$$

$$-\int (q \log q + q \log p) = -D(q||p) = -\int q(z, \theta) - \log \frac{q(z, \theta)}{p(z, \theta)} d\mu(z, \theta)$$

$$ELBO = \mathbb{E}_q \log p(x|z, \theta) - D(q||p)$$

Equivalent to minimizing

$$-ELBO = -\mathbb{E}_q \log p(x|z, \theta) + D(q||p)$$

Data term  
log probs of data

regularizer making variational distribution close to prior distribution

### Variational Autoencoders

Generative models that are trained via variational inference

"Decoder" is a neural net that generates from a latent variable

"Encoder" approximates posterior distribution with another neural network trained using variational inference

Example: Generative Model

$$z \sim N(0, I_n)$$

$$x|z \sim N(G(z), I_d)$$

generator network  
or  
decoder

$$G(z) = A_z \varphi(A_1 z + b_1) + b_2$$

We want to estimate  $p(z|x)$  but since  $G$  is non-linear, we estimate via variational inference

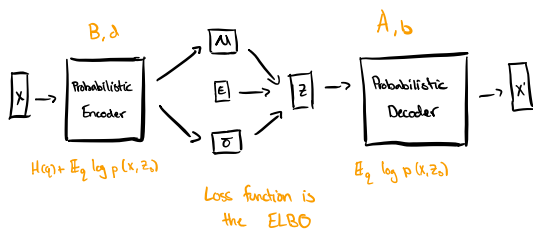
$$q(z|x) = N(\mu_x, \sigma_x^2 I_n)$$

variational parameters (shared across images)

We make an encoder neural network to output mean and variance

$$\mu(x) = B_x \varphi(B_1 x + d_1) + d_2$$

Similar for  $\log \sigma^2(x)$



We approximate  $\mathbb{E}_q(\log p(x, z))$  via sampling

$$\mathbb{E}_q(\log p(x, z)) \approx \frac{1}{N} \sum_{s=1}^N \log p(x, z_s)$$

parameters of network have disappeared!

$$\text{Reparameterize } z_i = \mu(x) + \sigma(x) \epsilon_i, \quad \epsilon_i \sim N(0, I_n)$$

Example: Suppose  $x|z \sim N(G(z), I)$  where

$$G(z) = \varphi(Az + b)$$

Then  $-\log p(x|z)$

$$\frac{1}{2} \|x - \varphi(Az + b)\|^2 + \text{constant}, \quad -\log(z) = \frac{1}{2} \|z\|^2 + \text{constant}$$

Further suppose  $q(x) = N(\mu(x), \sigma^2(x)I)$   
fit outputs

$$\mu(x) = \varphi(Bx + d)$$

$$\begin{aligned}
-\mathbb{E}_p \log p(x|z) &\approx \frac{1}{N} \sum_{s=1}^N \frac{1}{2} \|x - \varphi(Az_s + b)\|^2 \\
&= \frac{1}{N} \sum_{s=1}^N \frac{1}{2} \|x - \varphi(A(u(x) + \epsilon_s \sigma^2(x)) + b)\|^2 \\
&= \frac{1}{N} \sum_{s=1}^N \frac{1}{2} \|x - \varphi(A(\varphi(Bx + d) + \epsilon_s \sigma^2(x)) + b)\|^2
\end{aligned}$$

Regularizer

$$-\mathbb{E}_z \log p(z) = -\frac{1}{2} \mathbb{E}_z \|z\|^2 = \frac{1}{2} \|\lambda(x)\|^2 + \frac{\kappa}{2} \sigma^2(x)$$

$$-\mathbb{E}_z \log q(z) = \frac{\kappa}{2} \log \sigma^2(x)$$

entropy

Tensor flow then calculates the gradients for you

### Lecture 13: Sparsity and Graphs

graphs are a natural language for relationships

Undirected Graph

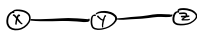
$$G = (V, E) \leftarrow V \text{ vertices and } E \text{ edges}$$

If  $x \in \mathbb{R}^p$ , we look at graphs with  $p$  vertices

Inverse of covariance matrix is the precision matrix

shows conditional independence relationships

Example:



$$X \perp\!\!\!\perp Z \mid Y$$

X and Z are independent conditioned on Y

Markov Property

Probability distribution  $P$  satisfies global Markov property with respect to a graph  $G$  if

For any disjoint vertex subsets  $A, B, C$  s.t.  $C$  separates  $A$  and  $B$

$$X_A \perp\!\!\!\perp X_B \mid X_C$$

$X_A$  are random variables  $X_j$  with  $j \in A$

$C$  separates  $A, B$  means every path between  $A$  and  $B$  passes through  $C$

Graph Estimation

Given  $n$  samples, estimate the graph  $G$

Gaussian Case

$$\Omega = \Sigma^{-1} \leftarrow \text{precision matrix}$$

0s in the precision matrix implies conditional independence

Parallel Lasso:

- For each  $j=1 \dots p$ , regress  $X_j$  on all other variables using lasso
  - Put an edge between  $X_i$  and  $X_j$  if each appears in the regression of the other
- intersection and union rule have the same asymptotic result

Graphical Lasso:

- Assume multivariate gaussian model
  - Subtract out sample mean  $\leftarrow$  center the data
  - Minimize negative log-likelihood of the data s.t. constraint on sum of absolute values of inverse covariance
- $\leftarrow$  encourages sparsity

Assuming a mean of 0

probability density  $\rightarrow$  
$$p(x) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp\left(-\frac{1}{2} \text{tr}(\Sigma^{-1} x x^T)\right)$$

$$-\log p(x) = \frac{1}{2} \log |\Sigma| + \frac{1}{2} \text{tr}(\Sigma^{-1} x x^T) = -\frac{1}{2} \log |\Omega| + \frac{1}{2} \text{tr}(\Omega x x^T)$$

using  $\log |A| = -\log |A^{-1}|$

Sum over all datapoints

$$\begin{aligned} -\sum_{i=1}^n \log p(x_i) &= \frac{1}{2} \sum_{i=1}^n \text{tr}(\Omega x_i x_i^T) - \frac{n}{2} \log |\Omega| \\ &= \frac{n}{2} \text{tr}(\Omega S_n) - \frac{n}{2} \log |\Omega| \end{aligned}$$

$S_n$  is sample Covariance

$$S_n = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$\leftarrow$  Convex function of  $\Omega$

Add a  $\lambda_1$  penalty for sparsity

Objective function  $\rightarrow$  
$$O(\Omega) = \text{tr}(\Omega S_n) - \log |\Omega| + \lambda \sum_{k,j} |\omega_{kj}|$$

rescaled by  $\frac{\lambda}{n}$

Parallel and graphical lasso produce different graphs

Choosing  $\lambda$

① Cross-Validation

tends to over select

② BIC = log-likelihood -  $(p/2) \log n$

③ AIC = log-likelihood -  $p$

## Lecture 14: Discrete Data Graphs and Graph Neural Networks

Hammersley, Clifford, Besag Theorem

A positive distribution over random variables  $z_1, \dots, z_p$  satisfies the Markov properties of graph  $G$  iff it can be represented as

$$p(z) \propto \prod_{c \in \mathcal{C}} \psi_c(z_c)$$

$\mathcal{C}$  is the set of cliques in the graph  $G$

$\uparrow$   
maximally connected  
subgraph

Discrete graphical models don't have a closed form  $\leftarrow$  no normalizing constant

Need to use Gibbs sampling and variational inference

Via Hammersley, Clifford and Besag Theorem

$$p(z; B) \propto \exp\left(\sum_{c \in \mathcal{C}} B_c \phi_c(z_c)\right)$$

Ising Model

$$p(z; B) = \exp\left(\sum_{i \in V} B_i z_i + \sum_{(i,j) \in E} B_{ij} z_i z_j\right)$$

Stochastic Approximation: Gibbs Sampling

Deterministic Approach: Mean field variational algorithm

No analog for the graphical lasso

Parallel Lasso: penalized logistic regression on each node  $z_i$  to estimate neighbors

- create edge  $(i,j)$  if  $j \in \hat{N}(i)$  and  $i \in \hat{N}(j)$  } Asymptotically equal

- (create edge (i,j) if  $j \in \hat{N}(i)$  or  $i \in \hat{N}(j)$ )

### Scaling Behavior

Ising Model:  $n \geq d^3 \log p$

Graphical Lasso:  $n \geq d^2 \log p$

Parallel Lasso:  $n \geq d \log p$

Lower bound:  $n \geq d \log p$

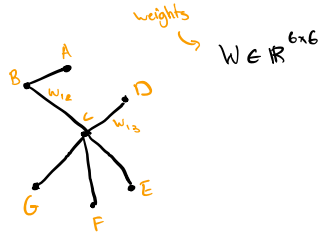
} Varying incoherence assumptions

### Graph Neural Networks

#### Graph Laplacian

$$L = D - W$$

degree matrix  $\rightarrow$   $D$   
weight matrix  $\leftarrow$   $W$



Weights  $\rightarrow W \in \mathbb{R}^{6 \times 6}$

	A	B	C	D	E	F	G
A	1	-1					
B	-1	2	-1				
C							
D		-1	1	1			
E					1	1	1
F					1	1	
G							1

For  $p$  nodes  $L \in \mathbb{R}^{p \times p}$

Polynomials of Laplacian are analogous to CNN filter

$$P_w(L) = w_0 I_n + w_1 L + w_2 L^2 + \dots + w_d L^d$$

weights  $w_i$  are analogous to filter coefficients

degree  $d$  of the polynomial is like the size of the kernel

Laplacian is a Mercer Kernel

$$f^T L f = \sum_{(u,v) \in E} (f_u - f_v)^2 \geq 0$$

#### Laplacian Operator

$$\frac{\partial^2}{\partial x^2} + \dots + \frac{\partial^2}{\partial x_p^2} \rightarrow \text{elegant connection}$$

#### Equivariance

A transformation is equivariant if

$$f(Px) = Pf(x)$$

where  $P$  is a permutation matrix  $PP^T = I$

#### Building Layers

$$h^{(k+1)} = \psi(P_w(L)h^{(k)})$$

## Lecture 15: Reinforcement Learning

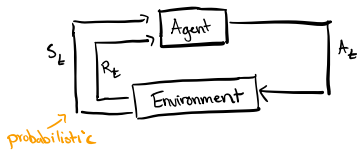
- Agent interacts with environment
- Actions the agent takes change the state of the environment
- Agent receives a reward for each action and seeks to maximize cumulative reward

Framework for sequential decision making to achieve a long term goal

#### Formalization of Reinforcement Learning

- Environment is in state  $s$  at a given time
  - Agent takes action  $a$
  - Environment transitions to state  $s' = \text{next}(s,a)$
  - Agent receives reward  $r = \text{reward}(s,a)$
- } Markov Decision Process

goal is to make optimal decisions from experience  
 In between supervised and unsupervised learning



Policy: mapping from states to actions

Reward Signal: Sequence of rewards received at each time step

Value function: A mapping from states to total reward

Values are predictions of future rewards

### Q-Learning

Maintains quality variable  $Q(s,a)$  for taking action  $a$  in states

Measure of cumulative rewards obtained by algorithm when it takes action  $a$  in states

### Algorithm

$$Q(s,a) \leftarrow Q(s,a) + \alpha (\text{reward}(s,a) + \gamma \max_{a'} Q(\text{next}(s,a), a') - Q(s,a))$$

"Gradient Ascent"

discount factor

"weighting of short term vs. long term rewards"

Bellman Equation: Deterministic Case

Value:

$$V_n(s) = \max_a \{ \text{reward}(s,a) + \gamma V_n(\text{next}(s,a)) \}$$

Fixed point equation

Q-function:

$$Q_*(s,a) = \text{reward}(s,a) + \gamma \max_{a'} Q_*(\text{next}(s,a), a')$$

Q-Learning algorithm satisfies value function bellman equation

### Lecture 16: Deep Reinforcement Learning

Q-Learning Algorithm is intended to estimate  $Q^*$

Parameters: step size  $\alpha$ , exploration probability  $\epsilon$ , discount factor  $\gamma$

Initialize  $Q(s,a)$  arbitrarily

Loop for each episode

Initialize  $S$

Loop for each step

Choose action  $A$  using  $Q$  with  $\epsilon$ -greedy policy

Take action  $A$ , observe reward  $R$  and new state  $S'$

$$Q(S,A) = Q(S,A) + \alpha (R + \gamma \max_{a'} Q(S',a) - Q(S,A))$$

$S \leftarrow S'$

Until  $S$  is terminal

↑  
best action for next state

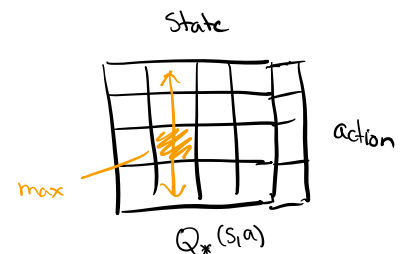
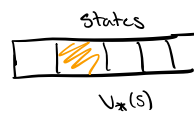
Q-learning algorithm is estimating bellman equation from data

If we know  $Q_*$ ,

$$V_*(s) = \max_a Q_*(s,a)$$

$$= \max_a \{ \text{reward}(s,a) + \gamma \max_{a'} Q_*(\text{next}(s,a), a') \}$$

$$= \max_a \{ \text{reward}(s,a) + \gamma V_*(\text{next}(s,a)) \}$$



Given  $Q_*$ , the system evolves deterministically

$$R_{total} = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$$

Bellman Equation: Random Environments

$$\begin{aligned} V_*(s) &= \max_a \sum_{s',r} p(s',r|s,a) \{r + \gamma V_*(s')\} \\ &= \max_a \mathbb{E} \left[ R_{t+1} + \gamma V_*(S_{t+1}) \mid S_t = s, A_t = a \right] \end{aligned}$$

Direct Application of Q-learning is only possible for small state and action spaces

Iterating over all possibilities

For large state space we map states to "features"

Deep Reinforcement learning uses a multilayer neural network to learn features and the Q-function

Goal: Represent  $Q_*$  as a neural network

$$Q(s,a; \theta_*) = \mathbb{E} \left[ R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a'; \theta_*) \mid S_t = s, A_t = a \right]$$

$\theta_*$  are weights in the neural network

Value iteration

i)  $Q \leftarrow \mathbb{E} \left[ R_{t+1} + \gamma \max_{a'} Q_{current}(S_{t+1}, a') \mid S_t = s, A_t = a \right]$

ii) Set  $Q_{current}$  to  $Q$  as updated in i)

Can be seen as stochastic gradient descent w.r.t least-squared loss

$$(Y(s,a; Q_{current}) - Q(s,a))^2$$

Strategy

target network:  $Q(\dots, \theta_{current}) \leftarrow$  used to evaluate bellman target

Policy network:  $Q(\dots, \theta)$

$\cdot Y_t = R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a'; \theta_{current})$

$\cdot$  Adjust parameters  $\theta$  to minimize squared error

$$(Y_t - Q(s,a,\theta))^2$$

Use SGD / backpropagation

$\cdot$  Update  $\theta_{current}$  to  $\theta$  and repeat procedure

Learning takes place when  $\mathbb{E}[Y] \neq Q$



# Lecture 17: Policy Gradient Method

A policy is a mapping from states to action

A value function is a mapping from states to total reward

Rewards are short term and values are predictions of future rewards

## Q-Learning

Parameters: step size  $\alpha$ , exploration probability  $\epsilon$ , discount factor  $\gamma$ , initialize  $Q(s,a)$  arbitrarily except terminal

Loop for each episode:

Initialize state  $s$

Loop for each step of episode:

Choose action  $a$  using  $Q$  with  $\epsilon$ -greedy policy

Take action  $a$ ; observe reward  $r$  and new state  $s'$

$$Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma \max_{a'} Q(s',a') - Q(s,a))$$

$$s \leftarrow s'$$

Until  $s$  is terminal

Q-Learning is an example of temporal-difference learning

"Off-policy" approach that is only practical if the space of states and actions are small

Deep Reinforcement learning uses a multilayer neural network to build a map from states to features for large state spaces

## General Strategy

$$Q(s,a;\theta) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) \mid S_t = s, A_t = a]$$

Sample  $y_t$  from the conditional distribution

$$y_t = R_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a'; \theta_{t+1})$$

Adjust parameters  $\theta$  to make squared error small

$$(y_t - Q(s,a;\theta))^2$$

↑  
via stochastic gradient descent

Learning happens when expectations are violated

Automatic Differentiation allows us to handle loss functions that are constructed dynamically

Calculate gradients on a "tape"

Use backpropagation on the "tape" to produce numeric gradients

## Policy Iteration

0. Initialize Policy arbitrarily

1. Compute values for current policy (Policy Evaluation)

Loop over all states and calculate expected reward

2. Update policy to match values (Policy Improvement)

Verify if policy agrees with Bellman equation

3. Return to step 1

## Policy Gradient Calculation

Episode:

$$(s_0, r_0) \rightarrow (s_1, r_1, a_1) \rightarrow (s_2, r_2, a_2) \rightarrow \dots \rightarrow (s_E, r_E, a_E) \rightarrow s_{E+1}$$

↑  
terminal

$$P(\tau|\theta) = \prod_{t=1}^T \underbrace{\pi(a_t|s_{t+1})}_{\text{policy}} \underbrace{P(s_{t+1}, r_{t+1} | s_t, a_t)}_{\text{environment}}$$

$$R(\tau) = \sum_{t=1}^T y_t$$

↗ reward

Objective function:

$$J(\theta) = \mathbb{E}[R(\tau)]$$

Use gradient descent to optimize objective function

$$\theta \leftarrow \theta + \eta \nabla_{\theta} J(\theta)$$

$$J(\theta) = \mathbb{E}[R(\tau)] = \int R(\tau) p(\tau|\theta) d\theta(\tau)$$

$$\nabla_{\theta} J(\theta) = \int R(\tau) \nabla_{\theta} p(\tau|\theta) d\theta(\tau)$$

$$= \int R(\tau) \frac{\nabla_{\theta} p(\tau|\theta)}{p(\tau|\theta)} p(\tau|\theta) d\theta(\tau)$$

$$= \mathbb{E}[R(\tau) \nabla_{\theta} \log p(\tau|\theta)] \leftarrow \text{can now sample trajectories}$$

$$\approx \frac{1}{N} \sum_{s=1}^N R(z^{(s)}) \nabla_{\theta} \log p(z^{(s)}|\theta)$$

$$= \frac{1}{N} \sum_{s=1}^N R(z^{(s)}) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

## Lecture 19: Reinforcement Learning: Actor-Critic Methods

### Multi-armed Bandits

- Rewards are independent and noisy
- Arm  $k$  has payoff  $\mu_k$  with variance  $\sigma_k^2$  on each pull
- Each step pull an arm and observe the reward
- Can estimate mean reward of each arm

Exploration - exploitation trade-off

### Policy Gradient Method: Loss Function

$$\text{Policy: } \pi_{\theta}(a_t | s_t)$$

$$\text{Environment: } p(s_{t+1}, r_{t+1} | s_t, a_t)$$

Episode:

$$(s_0, r_0) \rightarrow (s_1, r_1) \rightarrow \dots \rightarrow (s_T, r_T, a_T) \rightarrow s_{T+1}$$

terminal state

policy parameters

$$P(z|\theta) = \prod_{t=1}^T \pi_{\theta}(a_t | s_t) P(s_{t+1}, r_{t+1} | s_t, a_t)$$

Gradient Descent Over policy parameters

$$\nabla_{\theta} \log p(z|\theta) = \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Total Reward

$$R(z) = \sum_{t=1}^T r_t \quad (\text{no discounting})$$

Loss Function

$$J(\theta) = \mathbb{E}_{\theta}(R(z)) = \int R(z) p(z|\theta) d\theta$$

$$\nabla_{\theta} J_{\theta} = \int R(z) \nabla_{\theta} p(z|\theta) d\theta$$

$$= \int R(z) \frac{\nabla_{\theta} p(z|\theta)}{p(z|\theta)} p(z|\theta) d\theta$$

$$= \int R(z) \nabla_{\theta} \log p(z|\theta) p(z|\theta) d\theta$$

$$= \mathbb{E}_{\theta} \left( R(z) \nabla_{\theta} \log p(z|\theta) \right)$$

$$= \mathbb{E}_{\theta} \left[ R(z) \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

Sample

$$\approx \frac{1}{N} \sum_{s=1}^N R(z^{(s)}) \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(s)} | s_t^{(s)})$$

### Actor-Critic Approach

- Estimate policy and value function together

Actor is the policy used to select actions

Critic is the value function used to evaluate the actor

After a selected action, critic evaluates the state and updates actor and value function accordingly

Error Signal

$$d_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

$> 0$  better than expected

$< 0$  less than expected

Value Function update  
 $V(s_t) = V(s_{t-1}) + \alpha \delta_t$

Policy Parameter update

use gradient and  $\delta_t$

Lecture 19: Sequence Models and Recurrent Neural Networks

Language models are ways of assigning probability to any sequence of words

$$p(w_1, \dots, w_n) = p(w_1) p(w_2|w_1) \dots p(w_n|w_1, \dots, w_{n-1})$$

← conditional probability decomposition

Generate text by sampling conditional distribution

Perplexity is a rescaling of log-likelihood  
 geometric mean is no greater than the arithmetic mean

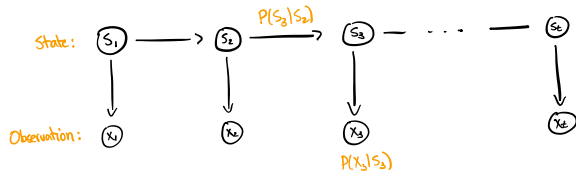
$$\text{Perplexity}(\theta) = \left( \prod_{i=1}^N p_\theta(w_i | w_{1:i-1}) \right)^{-1/N}$$

Inverse of the geometric mean of conditional probabilities evaluated on test data

Interpretation: if perplexity is 100, the model effectively predicts the next word with probability 1/100  
 geometric average "branching behavior" of model

Hidden Markov Models

- State is chosen stochastically
- Non-Markov on observations



probability of a word sequence

$$p(x_1 \dots x_n) = \sum_{s_1 \dots s_n} \prod_{t=1}^n p(s_t | s_{t-1}) p(x_t | s_t)$$

computed via forward-backward algorithm

Recurrent neural networks are similar to HMM but they are not stochastic

Lecture 20: RNNs, LSTMs, GRUs, and All That

Recurrent neural networks evolve states deterministically

Not mixture models like HMMs

State is distributed, not categorical like HMMs

$$\text{state: } h_t = \tanh(W_{hh} h_{t-1} + W_{hx} x_t + b_h)$$

3x3 matrix

4x3 3x1 previous state  
 3x4 4x1 current character  
 3x1 bias

$$y_t = W_{hy} h_t$$

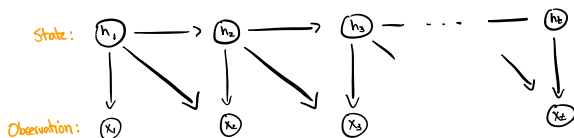
← vocabulary vector

$$x_{t+1} | y_t \sim \text{Multinomial}(\pi(y_t))$$

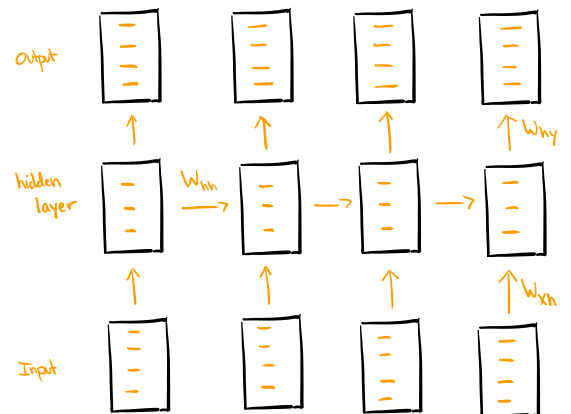
Soft-max

calculate derivatives w.r.t. probabilities and use for stochastic gradient descent

RNN Graphical Structure



Example RNN Structure



## LSTM: Long Short-Term Memory

gradients in traditional RNNs get really small or really big as the model learns

Adds a layer that "forgets" memory

## Gated Recurring Units (GRU)

High level idea:

- Learn when to update hidden state to "remember" important information
- Keep information until they are used
- Reset or forget info when it isn't useful

GRUs rely on gates

$\Gamma = 1$ : gate is open, information flows

$\Gamma = 0$ : gate is closed, information is blocked

Two-types of gates

$\Gamma^u$ : Long-term memory propagation when open, otherwise uses local state information

$\Gamma^r$ : When closed, local gate is reset. When open, the state is updated as a vanilla RNN

State evolution

$$c_t = \tanh(W_{hx} x_t + W_{hh} (\Gamma^r \odot h_{t-1}) + b_h)$$

candidate state

$$h_t = (1 - \Gamma_t^u) \odot c_t + \Gamma_t^u \odot h_{t-1}$$

determines

whether we use long term propagation

ignore previous state when closed

Hadamard Product

$$\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \odot \begin{pmatrix} 4 \\ 5 \\ 6 \end{pmatrix} = \begin{pmatrix} 4 \\ 10 \\ 18 \end{pmatrix}$$

Everything needs to be differentiable so each gate is calculated via a sigmoid function

$$\Gamma_t^u = \sigma(W_{ux} x_t + W_{uh} h_{t-1} + b_u)$$

$$\Gamma_t^r = \sigma(W_{rx} x_t + W_{rh} h_{t-1} + b_r)$$

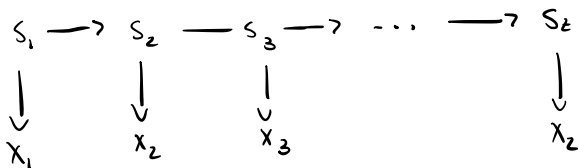
Sigmoid function

## Lecture 21: Sequence-to-Sequence models

Sequence models develop a conditional probability for the next state based on previous ones

## Kalman Filters

Gaussian version of HMM (generative model)



$$s_t | s_{t-1} \sim N(A s_{t-1}, \Gamma)$$

$$x_t | s_t \sim N(B s_t, \Sigma)$$

State is a real vector and evolves stochastically

Everything is linear

## Recurrent Neural Networks

- State is deterministic
  - Has an additional non-linearity
- } differentiates from Kalman Filter

## Gated Recurrent Units

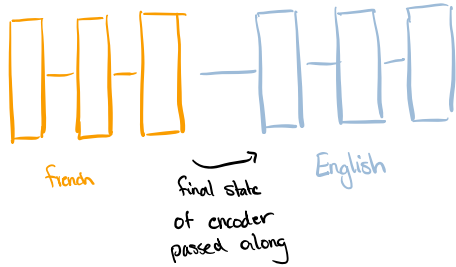
- Learn when to update hidden state to "remember" important information
- Keep information in memory until it is used
- Reset/forget information once it is no longer useful

$T=1$  : gate open ; information flows

$T=0$  : gate is closed ; information blocked

## Sequence-to-Sequence Models

Uses two RNNs : Encoder and Decoder

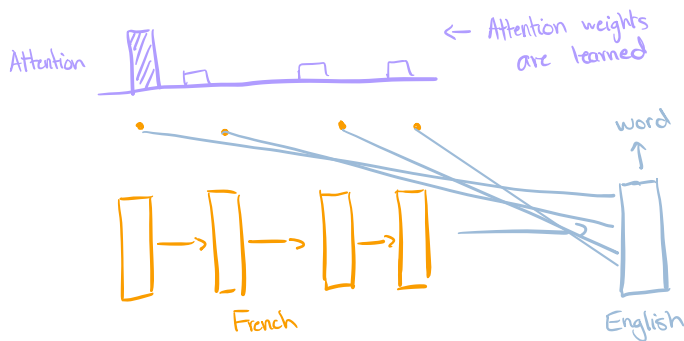


Decoder computes  $\prod_{t=1}^T p(y_t | h, y_1, \dots, y_{t-1})$

↑  
all information from French held in this state

## Attention/Alignment

Each step of decoding directly connects to encoder  
focuses attention on a certain segment



focuses more on the first word of the French sentence

## Attention Mechanism

- $m$  feature vectors or values  $V \in \mathbb{R}^{m \times v}$
- Model dynamically chooses which to use based on similarity between query  $q \in \mathbb{R}^q$  to a set of  $m$  keys  $K \in \mathbb{R}^{m \times k}$
- If  $q$  is most similar to  $i$ , then we use  $v_i$

# Kernel Regression as Attention

$$\hat{y}(x) = \frac{\sum_{i=1}^n k_h(x_i, x) y_i}{\sum_{i=1}^n k_h(x_i, x)} = \sum_{i=1}^n \alpha(x_i, x) y_i$$

attention weights
keys
query
values

In a gaussian kernel

$$\alpha(x_i, x) = \frac{e^{-\frac{1}{2k^2} \|x_i - x\|^2}}{\sum_j e^{-\frac{1}{2k^2} \|x_i - x\|^2}}$$

softmax

## Lecture 22: Attention and Transformers

Decoder state evolution

$$h_{t+1} = \tanh(W_{hh} h_t + W_{xh} x_t + W_{ah} a_t + b_h)$$

decoder
output
attention
bias

Attention mechanisms give a module for building new features

Algebraic representation of Attention

$$\text{Attn}(q, \{(k_i, v_i), \dots, (k_m, v_m)\}) = \sum_{i=1}^m w(q, k_i, m) v_i$$

$$w_i(q, k_i, m) = \frac{\exp(a(q, k_i))}{\sum_{j=1}^m \exp(a(q, k_j))}$$

Kernel Regression as attention

$$\hat{m}(x) = \sum_{i=1}^n w(x, x_i, n) y_i$$

query: test point x  
 keys: data  $x_1, \dots, x_n$   
 values: responses  $y_1, \dots, y_n$

If query and keys have the same dimension d

$$a(q, k) = q^T \frac{k}{\sqrt{d}} \in \mathbb{R}$$

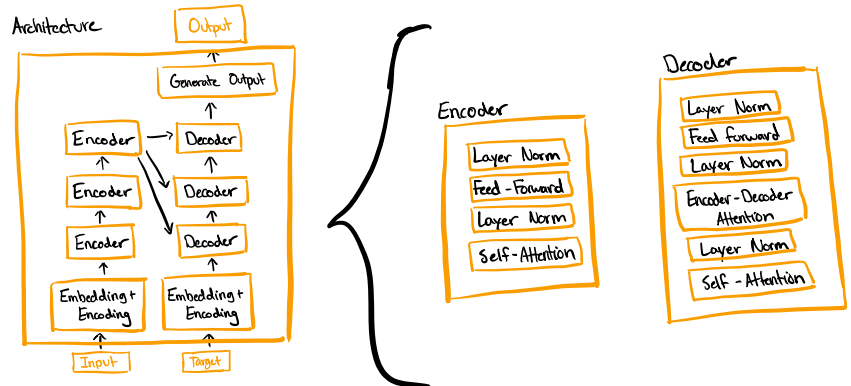
Matrix notation

$$\text{Attn}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$

weights

## Transformers

Transformers are RNN with no hidden structure



Transformers use self-attention instead of RNN  
multi-headed attention

## Embedding Layers

Embed words via vocabulary indices

Positional embedding is done via  $\sin/\cos$   
Neural networks don't do well with integers

For the  $t^{\text{th}}$  word in the sequence

$$PE_{(t, 2i)} = \sin\left(\frac{t}{L^{2i/d}}\right)$$

$d$  is the embedding dimension

$L$  is a maximum sequence length

$$PE_{(t, 2i+1)} = \cos\left(\frac{t}{L^{2i/d}}\right)$$

position is coded in the phase of the trig function

Concatenate embedding and positional encoding

## Multi-head Attention

Residual connections: add previous input to output of a step

Encodes multiple pairings of a sentence

## Encoder-Decoder Attention

### Lecture 23: Transformers and AI/ML Ethics

Attention Mechanisms give a flexible method for building new features

- $m$  feature vectors or values  $V \in \mathbb{R}^{m \times v}$
- Model dynamically chooses which to use
- Decide feature vector based on similarity between query vector  $q \in \mathbb{R}^q$  and set of  $m$  keys  $K \in \mathbb{R}^{m \times k}$
- If  $q$  is most similar to key  $i$ , then we use value/feature  $V_i$

Positional encoding via  $\sin/\cos$  functions

$$PE_{(t, 2i)} = \sin\left(\frac{t}{L^{2i/d}}\right)$$

$$PE_{(t, 2i+1)} = \cos\left(\frac{t}{L^{2i/d}}\right)$$

$d$  is embedding dimension

$L$  is maximum sequence length

transformer predicts sequentially

Residual connection brings back the input

Decoder masks future words

Attention all over the place in the input

Image vs. Caption Example

Self-attention: Image to Image Encoder  
Words vs. Words Decoder

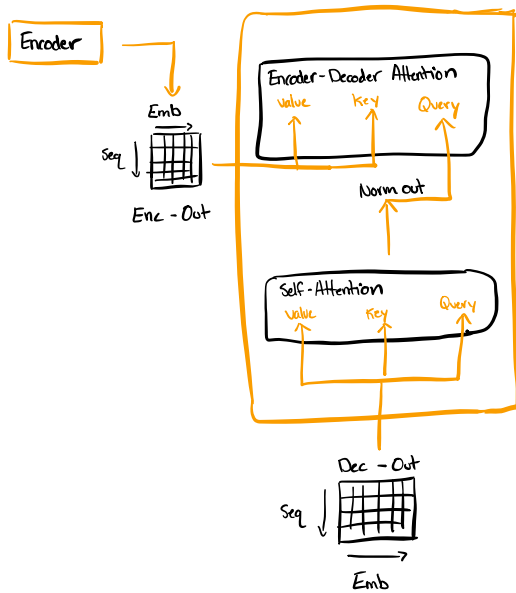
## Encoder-Decoder Attention: Image to Words

Differences between encoder and decoder

Decoder has temporal masking

Decoder conditions on encoder's state

## Encoder-Decoder Attention



Decoder - Out  $\rightarrow$  Linear  $\rightarrow$  Softmax  $\rightarrow$  Out

AI vs. ML

AI system exhibits a behavior resulting from decisions that are made

Embeddings encode societal biases

Lecture 24: Course Wrap Up