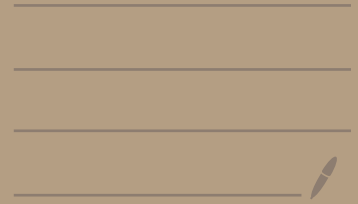


# STAT 457

---



# Week 1: K Nearest Neighbors and Bias/Variance tradeoff

## Core Regression Model

- $Y = f^*(X) + \epsilon$
- $X$  is feature vector
- $Y$  is numeric response
- $f^*$  some relationship
- $\epsilon$  is random noise

Identity:  $E[Y | X = x] = f^*(x)$

The regression function is the conditional mean.

## Training Data

- builds estimate for  $f^*$
- Used to predict new inputs

Training  $\rightarrow$  Estimator  $\rightarrow$  Prediction at  $x_0$

## K Nearest Neighbours (KNN)

- $N_k(x_0)$  - The  $k$  number of closest training points at  $x_0$
- Euclidean, assume no ties
- Prediction based on average nearby points

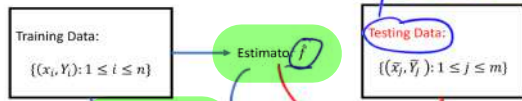
## Training error vs Testing Error

- The training error is evaluated on training data
- Always optimistic
- Testing error is far more accurate
- For KNN: Training error goes down as  $k$  goes down
- If  $k=1$  (1 group of cluster), training error = 0
- Testing error U shaped

Small  $k \rightarrow$  overfitting  $\rightarrow$  high test error

Large  $k \rightarrow$  underfitting  $\rightarrow$  high test error

## Training and testing errors



apply  $\hat{f}$  to make prediction for the data that is used to train  $\hat{f}$

$$\text{Training Error} = n^{-1} \sum_{i=1}^n (Y_i - \hat{f}(x_i))^2,$$

$$\text{Testing Error} = m^{-1} \sum_{j=1}^m (\tilde{Y}_j - \hat{f}(\tilde{x}_j))^2$$

► See code: knn\_regression.html. *Rk: testing error is what matters*

If  $k$  is chosen from training error. The training data is used twice. once to fit the model, then once to evaluate.

## Cross-Validation (Breaking loop)

- split training data into  $m$  folds
- for each fold train on  $m-1$ 
  - Test on remaining fold
- Average error

## Important Identity

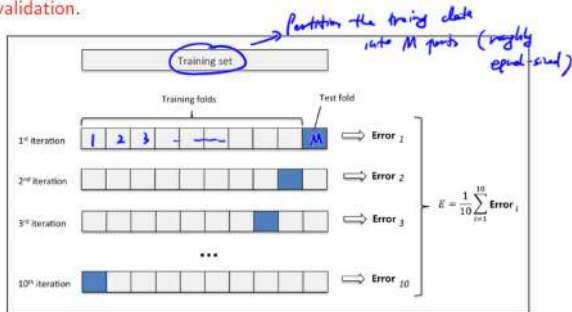
$$E[Y|X=x] = f^*(x)$$

$$Y = f^*(x) + \epsilon$$

Expectation is linear

True regression function is average output at a given input.

## validation.



Average out errors for  $m$  folds  
 Once you calculate all possible  $\alpha$ , cross validation candidates  
 Then choose lowest average CV error.  
 Then retrain all data.

In Sample Risk (Errin): Error on a special testing set that has same input but new  $Y_i$ .

The gap is the difference between In sample risk and training error.

Total Risk: gap + training error

### Quantifying Optimism

- $E[\text{Errin}] = E[\text{Err}_{\text{train}}] + \frac{2}{n} \sum \text{Cov}(Y_i, Y_i)$
- Covariance is how much models prediction changes based on training target ( $Y_i$ 's)
- High covariance means noise is followed very closely. (High optimism)

### Effective number of parameters ( $w$ )

- This is to minimize the total risk
- $w = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(Y_i, Y_i)$
- large  $w \rightarrow$  complex model, gap is large
- small  $w \rightarrow$  model is simple, gap is small

### Bias Variance Tradeoff

- 3 sources: irreducible error ( $\sigma^2$ )
  - $\hookrightarrow$  natural noise in dataset
  - $\rightarrow$  Bias: too simple does not capture patterns
  - $\rightarrow$  variance: too sensitive,

### Empirical Risk Minimization

- low training error does not always mean good
- Simple model  $F$  (high training error)
- Rich/complex models (overfitting)

# Week 2: Linear Algebra Review

## Trace

- For square matrix, SUM of its diagonal entries

$$\text{Tr}(AB) = \text{Tr}(BA)$$

## Linear Representation

- Vectors can be linear combinations of each other

$$v_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, v_2 = \begin{bmatrix} 4 \\ 1 \\ 5 \end{bmatrix}, v_3 = \begin{bmatrix} 2 \\ -3 \\ -1 \end{bmatrix}$$

$$v_3 = -2v_1 + v_2$$

- Vector is linearly indep, if it can't be made up using other vectors.
- Rank is number of lin indep vectors in matrix

$$\text{Col Rank}(A) = \text{Row Rank}(A) = \text{Rank}(A)$$

## Subspaces

- two vectors are closed under addition and multiplication
- $x+y \in V$ ,  $\alpha x \in V$
- dim of subspace  $V$  is largest subset of lin indep vectors
- **Span**: Smallest subspace contains every possible linear combination of those vectors.

**Basis**: span is equal to subspace. Minimum building blocks to create every vector in that subspace.

## Orthogonality

- two vectors are orthogonal if their inner product  $x^T y = 0$ . A vector is orthogonal to a subspace where its perpendicular to every vector inside the subspace.

Orthonormal matrix: every vector perpendicular to every other vector and unit magnitude of one. (Preserves angle and length)

Property:  $U^T U = U U^T = I$

- For 2 by 2 orthonormal matrices, these are ones that either rotate or reflect
- Preserves length  $\|Ux\|_2 = \|x\|_2$
- Preserves angle  $\langle Ux, Uy \rangle = \langle x, y \rangle$

## SVD

- Breaking down matrix into  $X = U \Sigma V^T$
- $U$  has orthonormal columns,  $V$  is an orthonormal matrix,  $\Sigma$  is diagonal matrix.

Spectral Decomp: symmetric square matrix,  $Ax = \lambda x$  ( $A$  times eigenvector)

Positive Def:  $x^T A x$  is always  $> 0$

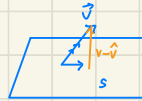
Positive Semi-Def:  $x^T A x$  is always  $\geq 0$ .

Eigen values also determine definiteness

## Projections

- vector  $v$ , subspace  $S$

$$v - \hat{v} \perp S$$



The residual is perpendicular to subspace, because removes all components of  $S$ .  
Minimizes distance

$$\|v - u\|^2 = \|v - \hat{v}\|^2 + \|\hat{v} - u\|^2$$

$$\|v - u\| \geq \|v - \hat{v}\|$$

Iterated Projection: Drop  $v_1$  into  $S_2$ , then drop result in  $S_1$ .

# Week 2/3: Optimization

- ERM (Empirical Risk Minimization)
  - ↳ finds prediction rule  $f = \xi_{f_0: \theta \in \Theta}$  that minimizes average loss.
- gradient: direction and steepness, vector of first derivatives
- Hessian  $H$ : tells you curvature

## Taylor Expansions

- approximate complex functions
- fastest descent ( $-\nabla F$ )
- $f(B + \delta) \approx f(B) + \nabla f(B)^T \delta$
- $f(B + \delta) \approx f(B) + \nabla f(B)^T \delta + \frac{1}{2} \delta^T H f(B) \delta$
- $g(\delta) = f(B) + \nabla f(B)^T \delta$

$$\nabla f(B)^T \delta \text{ subject to } \|\delta\| = 1$$

$$\delta^* = - \frac{\nabla f(B)}{\|\nabla f(B)\|}$$

$$\delta = -\nabla f(B)$$

$$\delta \propto \nabla f(B)$$

is proportional soon

## Local minima

$$f(B) \geq f(B_0) \text{ if } \|B - B_0\| < \delta$$

If you pick any  $B$  close enough to  $B_0$ .

the function value at  $B$  is not smaller than  $B_0$

local minima  $\cong$  global minima

(maybe)

Necessary condition:  $\nabla f(B_0) = 0$ , slope is 0, can't have immediate decrease

Sufficient condition:  $\nabla f(B_0) = 0$  and  $H f(B_0)$  is positive definite,  $B_0$  is local min  
(guarantee)

## Gradient Descent Algorithms

take small steps down till flat hill

First order Algorithm:  $x_{n+1} = x_n - \alpha_n \nabla f(x_n)$

Stop when  $\|\nabla f(x_n)\| \in \epsilon$

next guess      current guess      step size      grad at current point

## Newton Raphson

uses second order

$$x_{n+1} = x_n - H f(x_n)^{-1} \nabla f(x_n)$$

Hessian curvature

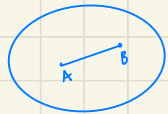
uses slope and curvature, faster but more expensive

# Convex Sets

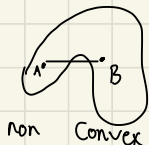
• A set  $A \subset \mathbb{R}^d$  is convex if:

for any two points  $x, y$  in the set, every point on the straight line between them is also in the set

$$\alpha x + (1-\alpha)y \in A, \text{ for } \alpha \in (0,1)$$



Convex



non Convex

$A_1 \cap A_2$  is convex  
if they are individually convex.

• Convex functions:  $f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$

value in middle

straight line interpolation

• Convex means: curve stays under line

• Examples:  $f(x) = |x|$   
 $f(x) = \max(x, 0)$

• Convex set  $\{x : f(x) \leq c\}$

## Proving Convexity

1) Hessian is PSD

2) First order tangent line condition

$$f(y) \geq f(x) + \nabla f(x)^T (y-x)$$

if  $\nabla f(x_0) = 0$ ,  $x_0$  is global min

# Week 3: Linear Regression

• linear regression  $f_{\beta}(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

$\beta = (\beta_0, \beta_1, \dots, \beta_p)$  unknown parameters

• model predicts number using weighted sum

• Loss function (RSS):  $RSS(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2$

• Matrix formulation:

$$X = \begin{bmatrix} 1 & x_{11} \\ \vdots & \vdots \\ 1 & x_{n1} \end{bmatrix} \quad (n \times (p+1))$$

• response vector:

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Prediction:  $X\beta =$  all predicted values

$$RSS(\beta) = \|Y - X\beta\|^2$$

• RSS is a quadratic function of  $\beta$ .

$$\hookrightarrow \text{gradient } \nabla_{\beta} RSS(\beta) = -2X^T Y + 2X^T X \beta$$

$$H_{RSS}(\beta) = 2X^T X$$

• Solution is unique when  $\text{rank}(X) = p+1$ , columns of  $X$  are lin indep

• Projection matrix:  $P = X(X^T X)^{-1} X^T$   
 $Y = PY$

## Subset Selection

- Choosing subset of features instead of all of them in lin reg
- reduces error, avoid overfitting

### 2 ways

1. Search strategy (try diff subsets of predictors)
2. Scoring rule (find best subset)

• Score = Fit + Penalty

$RSS(\hat{\beta}_m) =$  residual sum of squares for  $m$

• Common choices include: Mallows CP, AIC, BIC

## KNN

7/9

- non Parametric, Supervised, no functional form
- K is only Key Param
- distance based Voting / averaging
- No training Phase (stores data)

## Cross fold

- Split into  $m$  folds
- train on  $m-1$ , test on leftover
- average out errors, choose lowest  $K$

L.R - has effective number of params -  $p+1$

KNN - effective number of Params depends on  $K$

$$E.P = \frac{p}{K}$$

↳ how flexible model is

small  $K$  - low bias, high variance (wiggly)

Big  $K$  - high bias, low variance (smooth)

Bias = error from simplifying assumptions

Var = sensitivity to training data.

$Ax$  = matrix · scalar (weight)

$x^T y = 0$  orthogonal

$U^T U = U U^T = I \rightarrow$  orthonormal preserves angle/length

PSD  $\Rightarrow$  convex function

$P \rightarrow$  projection matrix, closest point to subspace

$a^T x + b \rightarrow$  affine  $\rightarrow$  convex

---

Gradient Desc  $x_{k+1} = x_k + \alpha \nabla f(x_k)$

Newton Raphson  $x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$

If convex, global min = local min

# Midterm 2

---



# Quiz Content (7 bullets)

## Linear Regression

- supervised learning predicts numerical response

model formulation:  $y = X\beta + \epsilon$

$y \in \mathbb{R}^n$ : response vector

$x \in \mathbb{R}^{n \times (p+1)}$ : design matrix

$\beta \in \mathbb{R}^{p+1}$ : coefficients

$\epsilon$  is noise

• minimize:  $RSS(\beta) = \|y - X\beta\|^2$

• choose Best  $\beta$  parameters

L.R: Objective:

Proof for solving:

Trying to choose Best  $\beta$  coefficients so that we minimize RSS

## Projection Matrix

$y = X\beta$

sub back in:  $y = X(X^T X)^{-1} X^T y$

$H = X(X^T X)^{-1} X^T y$

$y = Hy$

Projection matrix projects  $y$  onto column space of  $x$

## Ordinary Least Squares Solution

$$\nabla RSS(\beta) = -2X^T(Y - X\beta)$$

$$X^T(Y - X\beta) = 0$$

$$X^T Y - X^T X \beta = 0$$

$$X^T Y = X^T X \beta$$

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

$$RSS(\beta) = \|y - X\beta\|^2$$

$$-2X^T(Y - X\beta)$$

$$X^T(Y - X\beta) = 0$$

$$X^T Y - X^T X \beta = 0$$

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

L.R:  $\|y - X\beta\|^2$

R.R:  $\|y - X\beta\|^2 + \lambda \|\beta\|^2$

L.R:  $\|y - X\beta\|^2 + \lambda \|\beta\|^2$

## Linear regression

$y = X\beta + \epsilon$

Minimizing RSS:

$\min \|y - X\beta\|^2$

OLS:  $\beta = X^T(X^T X)^{-1} y$

projection matrix:

$H = Hy$

where  $H = X^T X (X^T X)^{-1}$

Projection matrix is symmetric and independent

$H^2 = H, H^T = H$

$\text{trace}(H) = p$

number of Params

high variance, low bias

## Ridge Regression

$y = \beta x + \lambda \|\beta\|^2$

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

as  $\lambda$  increases, less flexible, less effective param

Harsher penalty, lower variance

shrinks by  $\sum \frac{d_j}{d_j + \lambda}$

For  $\lambda = 0, df = p$

# Ridge Regression

OLS:  $\min_{\beta} \|y - X\beta\|^2$

Ridge Adds: L2 Penalty

$$\hat{\beta}_x = \arg \min (\|y - X\beta\|^2 + \lambda \|\beta\|^2)$$

$$\lambda \geq 0$$

$$\|\beta\|^2 = \beta^T \beta$$

Penalize large coefficients to reduce variance  
shrinks them toward zero

given  $B^{OLS} = (3, 1, 5)$   $\lambda = 2$

Show variable selection

$$B_{i1} = +1 (+1) = 1$$

$$B_{i2} = (+1)(-1) = -1 \rightarrow 0$$

$$B_{i3} = (+1)(+1) = 1$$

$S\lambda$  = matrix that maps  $y$  to fitted range

## Closed Form Solution

$$\nabla B = -2x^T(y - x\beta) + 2\lambda\beta$$

$$0 = -2x^T(y - x\beta) + 2\lambda\beta$$

$$x^T x \beta + \lambda \beta = x^T y$$

Factor  $(x^T x + \lambda I)\beta = x^T y$

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

For ridge,  $y\hat{z} = S\lambda y$

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

Not orthogonal projection

Ridge shrinks regularization projection toward 0

Model behaves as if it has less effective parameters. Shrinkage

$$dK/d\lambda = \sum_{j=1}^p \frac{d_j}{d_j + \lambda}$$

## Ridge Main Takeaways

- shrinks coefficients
- smaller coefficients  $\rightarrow$  lower variance
- shrinkage pulls away from truth  $\rightarrow$  bias

Higher  $\lambda$  = more Penalty  
-lower variance, higher bias

$$\arg \min \|y - X\beta\|^2 + \lambda \|\beta\|^2$$

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

$$x^T x + \lambda I = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0 \\ 0 & 2 \end{bmatrix}$$

inverse

$$\begin{bmatrix} \frac{1}{5} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \cdot \begin{bmatrix} 8 \\ 2 \end{bmatrix} = \begin{bmatrix} \frac{8}{5} \\ 1 \end{bmatrix}$$

## Two Forms

1.) Penalized:  $\min_{\beta} \|y - X\beta\|^2 + \lambda \|\beta\|^2$

2.) Constrained:  $\min_{\beta} \|y - X\beta\|^2$  s.t.  $\|\beta\|^2 \leq t$

only allowed to choose  $\beta$  inside ball of radius  $\sqrt{t}$

for every  $\lambda$ , there exists radius  $t$

only penalize coeff, not intercept, that would only shift up/down

eigenvals. of  $X^T X$ :  $\sum \frac{d_j}{d_j + \lambda}$

singular vals. of  $X$ :  $\sum \frac{\sigma_j^2}{\sigma_j^2 + \lambda}$

more intuitive

# Lasso Regression

- stands for least absolute shrinkage selection operator

## Core Definition

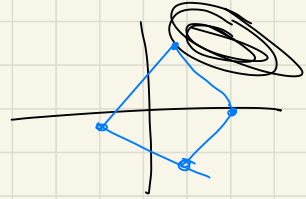
- minimizing RSS with an  $L_1$  norm penalty on the coefficients making the objective function is: 
$$\min_{\beta \in \mathbb{R}^p} \|y - X\beta\|^2 + \lambda \|\beta\|_1$$

$\lambda$  is a tuning parameter, controls bias/variance trade off. Larger  $\lambda$ , more shrinkage

## Key Properties

- Variable Selection: unlike ridge, some coefficients can be set to 0  
↳ Performs automatic feature selection. if  $\hat{\beta}_j = 0$ , feature does not affect final predictions
- Lasso constraint region is diamond shape. geometry?
- Lasso objective is convex because RSS term and  $L_1$  are convex.  
Meaning: local min found is global minimum

- Equivalent Constrained Form:



## Solution under orthogonal design

- Closed form solution exists if design matrix is orthogonal ( $X^T X = I$ )  
↳ Solution for each coefficient given by soft-thresholding

Each Prediction is uncorrelated and contains completely independent information

$$\hat{\beta}_j^{\text{lasso}} = \text{sign}(\hat{\beta}_j^{\text{OLS}}) (|\hat{\beta}_j^{\text{OLS}}| - \tau/2)$$

if magnitude of least squares is smaller than  $\tau/2$ , lasso estimate becomes 0

Ex. given  $X^T X = I$   
 $\hat{\beta}_{\text{OLS}} = (3, -1, 0.5)$   
 $\tau = 2$

$$\begin{aligned} \hat{\beta}^{\text{lasso}} &= +1 (2) = 2 \\ \hat{\beta}^{\text{lasso}} &= -1 (0) = 0 \\ \hat{\beta}^{\text{lasso}} &= -1 (-0.5) = -0.5 < 0 = 0 \end{aligned}$$

- Lasso lacks closed form solution, therefore uses coordinate descent (good for large data)

converts High Dimensional  $\rightarrow$  One dimensional Problems

- Optimize one coefficient at a time

# Coordinate Descent Process

1) Compute residual:  $r = y - X\beta$

$$\beta_j = \text{sign}(z_j)(|z_j| - \lambda)$$

2) Measure correlation with residual:  $z_j = \frac{1}{n} x_j^T r$

$z_j$  is partial residual correlation

Big value  $\rightarrow$  feature matters  
Small val  $\rightarrow$  feature useless

$$r = y - X\beta$$
$$z = X^T r$$

3) Apply soft thresholding

$$\beta_j = \begin{cases} z_j - \lambda & , z_j > \lambda \\ 0 & |z_j| \leq \lambda \\ z_j + \lambda & z_j < -\lambda \end{cases}$$

where  $z_j$  is correlation

$$\beta_{OLS} = (\text{sign } \beta_j)(|\beta_j| - \lambda)$$
$$= + (3 - 1)$$
$$= 2$$
$$\left. \begin{array}{l} 0 - 3 - 1 \\ = -0.5 \\ = 0 \end{array} \right\} \rightarrow (2, 0)$$

## Gradient Descent

• update all parameters at once

$$\beta_{k+1} = \beta_k - \alpha \nabla f$$

next step      last step      - learning rate      times gradient

$$r = y - X\beta$$
$$z = X^T r$$
$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 0 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 \end{bmatrix}$$

$$r = \begin{bmatrix} 3 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 & 0 \end{bmatrix}$$
$$r = \begin{bmatrix} 3 \\ 0 \end{bmatrix}$$
$$\begin{bmatrix} 3 & 1 \\ 2 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 3 & 3 \\ 2 & 4 \end{bmatrix}$$

$$z = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 0 \end{bmatrix}$$

$$z = \begin{bmatrix} 3 & 3 \\ 2 & 0 \end{bmatrix} \rightarrow \text{soft threshold} \\ (1, 0)$$

# Polynomial Regression - type of linear regression

- Linear regression we assume data can be fit using straight line  $y = mx + b$
- When data is non-linear or curves, we use a bending term  $x^2, x^3, \dots$

Higher Degree = more flexible But can overfit.

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$$

- uses expanded features, that are transformations of  $x$   
↳ Basis expansion

- Splines are broken down parts of small polynomials together  
↳ Breaks  $x$  axis into pieces.

- The joining points of these polynomials is knots

- Cubic splines: each piece is knot of degree 3 polynomial

↳ smooth curvature:

function is continuous on both sides of knot

first derivs continuous

second derivs continuous

- Spline basis functions: pre built building blocks automatically create piecewise cubic behaviour and enforce smoothness at knots.

Natural cubic spline: forced to become linear at boundaries

→ piecewise cubic polynomials

→ smoothness at knots

→ boundary constraint, second deriv = 0

## How to find "solution"

- 1.) choose  $k$  knots
- 2.) build spline basis function
- 3.) fit using OLS:  $\hat{\beta} = (X^T X)^{-1} X^T y$

Dimension of Cubic Spline Space

$$\bullet k + 4 \rightarrow \underbrace{(4k + 1)}_{\text{Params Knots}} - \underbrace{3k}_{\text{3 conditions}}$$

- Natural cubic:  $k + 2$   
↳  $k + 4 - 2$  } 2 new boundaries  
second deriv = 0 on both sides

# Smoothing Splines

• Fit data, penalize curve

$$\min \underbrace{\sum (y_i - f(x_i))^2}_{\text{fit data}} + \underbrace{\lambda \int f''(x)^2 dx}_{\text{smoothness penalty}}$$

• Degrees of freedom for Smoothing Splines, same as ridge  $\hat{g} = S_\lambda y$ ,  $df = \text{tr}(S_\lambda)$

•  $\lambda$  Controls tradeoff

$\lambda$  large  $\rightarrow$  almost linear  $\rightarrow$  lower variance

$\lambda$  small  $\rightarrow$  high curvature  $\rightarrow$  higher variance

• knots at every data point.

• Controls Complexity by Penalizing Curvature from  $\lambda$ .

## Questions - Practise

Ridge regression: Objective function is:  $J(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \sum \beta_i x_{ij})^2 + \lambda \sum \beta_j^2$

$\lambda = 0 \rightarrow$  same as OLS

$\lambda$  increase means less flexible, effective number of params decrease  
Coeff's shrink toward zero

4)  $f(\alpha, \beta) = (1 - \alpha - \beta)^2 + |\alpha| + \beta^2$

Using coordinate Descent:

Fix  $\alpha$ , minimize over  $\beta$

$$f(\alpha) = (1 - \alpha - \beta)^2 + |\alpha| + \text{const}$$

Minimize:  $(1 - \alpha - \beta)^2 + \beta^2$   
 $c = 1 - \alpha$   
 $(c - \beta)^2 + \beta^2$

Expanding, we get  $\beta = \frac{1 - \alpha}{2}$

Plug back into  $f$ :  $1 - \alpha - \beta = 1 - \alpha - \frac{1 - \alpha}{2}$

∴ long expansion:  $\min(\alpha - 1)^2 + 1 \cdot |\alpha|$

$1 - 1 = 0$   
 $\alpha^* = 0 \therefore (0, 1/2)$

3.) 10 knots, 11 intervals

$$g(x) = \beta_i x + \alpha_i \text{ means}$$

1 slope param  $\beta_i$

1 intercept param  $\alpha_i$

$$2 \times 11 = 22$$

$$22 - 10 \text{ knots}$$

$$\text{Dimension} = 12$$

$m$  knots, 2 extra conditions on front and last, meaning

$$\text{dim} = m$$

# Regression Trees

- Building predictive models by recursively partitioning feature space into smaller regions
- Top-down greedy approach
- Partition feature space into non overlapping sub regions.  $R_1, R_2, \dots, R_k$
- Tries to choose smart splitting rule to give greatest reduction in RSS
- Greedy nature is from making best split at that moment.
- Predicted value  $\hat{c}_k$  is average of response variables ( $y_i$ )
- Growing tree until every leaf is small leads to overfitting (high variance)
  - ↳ Cost complexity pruning

Complexity pruning: instead of minimizing RSS, minimize penalized version  
 $PRSS(T) = RSS(T) + \alpha |T|$ , where  $|T|$  is number of leaf nodes and  $\alpha$  is tuning parameter

$\alpha$  Large values  $\rightarrow$  Smaller trees

$\alpha$  small values  $\rightarrow$  Larger trees

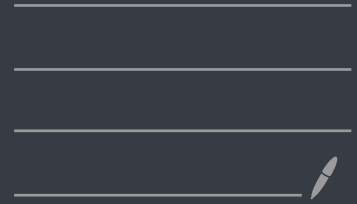
## Key Properties

- Very interpretable and easier to explain
- automatic feature handling
- invariance
- high variance
- Predictive power

• Effective number of  
Params = number of  
leaf nodes.

# Midterm 3

---



# PCA - Principal Component Analysis

- Finds new principal components that capture as much variance as possible
- Reduces dimensionality, clean directions with most signal
- All principal components are mutually perpendicular
- Works best for linear models / distance based

## Computation Process for PCA

- Center the data (Subtract mean from features)
  - Compute covariance matrix:  $\rightarrow$
  - Find eigenvectors and eigenvalues of covariance matrix
- $\downarrow$  Principal Components
- $\downarrow$  how important that PC is (variance)

$$\Sigma = \frac{1}{n} X^T X$$

$n$  is number of observations (rows)

Explained Variance:  $\frac{\lambda_i}{\sum_j \lambda_j}$

- Sort by importance, then project data onto top  $K$  PCs:  $\rightarrow z = X V_K$

$V_K$  is top  $K$  eigenvectors

Ex

Given  $\lambda_1 = 5$ ,  $\lambda_2 = 2$ ,  $\lambda_3 = 1$

Find amount of variance per PC

Find number of components needed to keep at least 85% of var

i)  $PC1: \frac{5}{8} = 62.5\%$   
 $PC2: \frac{2}{8} = 25\%$   
 $PC3: \frac{1}{8} = 12.5\%$

}  $62.5 + 25 = 87.5\%$   
Keep PC1 and PC2

To find eigenvalues  $\det(A - \lambda I) = 0$   
 $(A - \lambda I) V = 0$

Variance Lossed = Sum of unused eigenvalues

# Need to Memorize

## Linear Regression

- Model:  $y = XB + \epsilon$
- Prediction  $\hat{y} = X(X^T X)^{-1} X^T y$
- Hat matrix  $H = X(X^T X)^{-1} X^T$
- Hat matrix projects  $y$  onto column space of  $X$
- Effective number of param =  $\text{trace}(H) = p$ 
  - ↳ Equal to number of  $p$  features
- 2 main properties of Hat matrix.  $H = H^T$ ,  $H^2 = H$
- affine function:  $\beta_0 + \beta_1 z_1 + \dots + \beta_p z_p$  often better in practice

- Overfitting = higher variance / low bias  
- Underfitting = high bias, low variance

- Objective:  $\min \|y - XB\|^2$

- Closed form solution:  $\beta = (X^T X)^{-1} X^T y$