

Search David Fleet to get to website.

[www.cs.toronto.edu/~fleet/courses/c11/index.html](http://www.cs.toronto.edu/~fleet/courses/c11/index.html)

Learn the notes!

Office hour: 2-3 Friday

3 assignment:

- 3 weeks to work on it  $\approx 36\%$
- Quiz after assignment
- Midterm: 15%
- Final: 49%

piazza site is set up for this class.

### 3 main types of learning

$$y = f(x; \theta)$$

$\uparrow$        $\downarrow$   
 output    input

$\theta$ : set of parameters  
 $y$ : outcome you care about

data  $\{(x_i, y_i)\}_{i=1}^N$   $N$  - training examples

### Two classes in supervised learning

① classification - where  $y$  is one of distinct number of classes  
 $y \in \{1, \dots, C\}$  e.g. Binary classification, is this email spam or not spam

② regression - where  $y$  lives in real value vector space and contains continuous variables  
 $y \in \mathbb{R}$ . Eg acceleration in car not only 2 speeds.

## 2<sup>nd</sup> type of learning: Unsupervised Learning

No target data.

$$\{x_i\}_{i=1}^N$$

Ex/ Bunch of images or text documents is given to you.

Discovering with data when you don't have a target

3 things with unsupervised data

### 1) Clustering

- grouping or finding clusters that are similar

### 3) Density Estimation

This sequence is more probable than this.

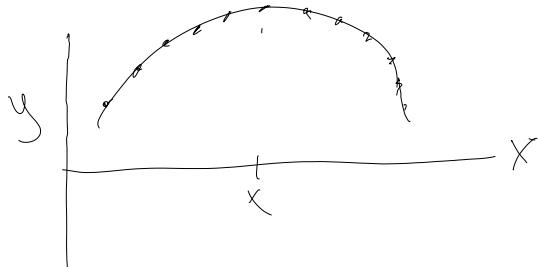
### 2) feature selection / Dimensionality Reduction

finding low dimensional features of high dimensional data in compression.

JPEG is a form of dimensionality reduction

e.g. this is an unlikely image in natural. But this image is highly probable.

### Model

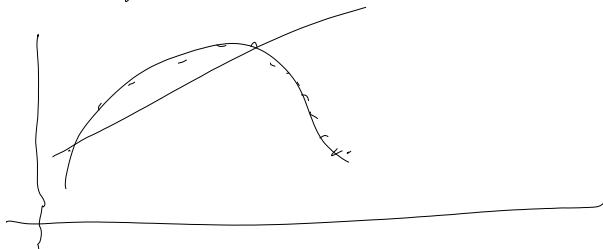


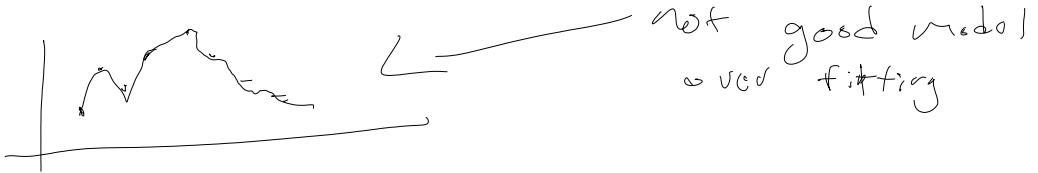
e.g.

A family of functions used to take parameters to make a prediction.

(Good model)

Decision, should I fit a line? or curve?





under fitting: errors in training data is too large

over fitting: fit the data too well.

It's not about fitting accurate model but predicting unseen test data.

In other words don't fit noise.

Eg fit polynomial model



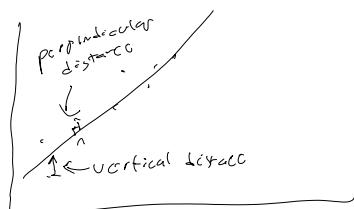
### Loss function

measures and quantifies performance on data set

$$y = f(x; \theta)$$

parameter = one that minimizes loss function  
calculated error over the training set.

loss funct ex/



Ex Bin prob  $y \in \{-1, 1\}$

$$\text{sgn}(yf(x))$$

1 // . . . . . -1 .. 1

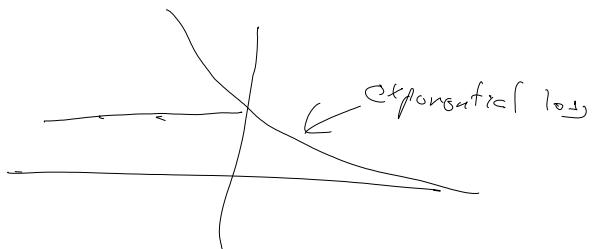
$\forall x \in \mathbb{R}^m$  given  $y \in C^{-1}, 1)$



$$\operatorname{sgn}(y f(x))$$

$$\frac{1}{2}((1 - \operatorname{sgn}(y f(x)))^2)$$

loss function, # of curves you created



problem: not differentiable.

how to get a good loss function?

IDK figure it out later. But its important.

### Linear Regression

$$\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}, y_i \in \mathbb{R}$$

model:  $y = f(x) = w x + b$

goal: find good values for  $w$  &  $b$ .

define error

$e_i = \text{training error}$

$$e_i = y_i - f(x_i) = y_i - (w x_i + b)$$

loss: sum of the squared error

$$E(w, b) = \sum_{i=1}^N e_i^2 = \sum (y_i - (w x_i + b))^2$$

$$\nabla_{w,b} E = 0$$

$$\left( \frac{\partial E}{\partial w}, \frac{\partial E}{\partial b} \right)$$

$$= \sum_{i=1}^N (y_i^2 - 2y_i(wx_i + b) + (wx_i + b)^2)$$

$$= \sum_{i=1}^N y_i^2 - 2y_iwx_i - 2y_ib + w^2x_i^2 + 2bwx_i + b^2$$

$$\frac{\partial E}{\partial b} = \sum_i 0 + 0 - 2y_i + 0 + 2wx_i + 2b$$

$$= -2 \sum_i (y_i - wx_i - b)$$

$$= -2(\bar{y}_i - w\bar{x}_i - Nb)$$

$$b^* = \frac{1}{N} (\bar{y}_i - w\bar{x}_i)$$

$$\bar{y} = \frac{1}{N} \sum y_i \quad b^* = \bar{y} - w\bar{x}$$

$$x = \frac{1}{N} \sum x_i$$

objectivo

$$(E(w, b)) = \sum_{i=1}^N e_i^2 = \sum (y_i - (wx_i + b))^2$$

$$E(w, b^*) = \sum_i (y_i - wx_i - \bar{y} + w\bar{x})^2$$

$$= \sum_i \{(y_i - \bar{y}) - w(x_i - \bar{x})\}^2 \quad \text{Ex/ expand}$$

$$\frac{\partial E}{\partial w} = -2 \sum_i ((y_i - \bar{y}) - w(x_i - \bar{x})) (x_i - \bar{x})$$

$\frac{\partial}{\partial w} \sum_i f(x_i, w)^2$ 

$$= 2 \sum_i f(x_i, w) \frac{\partial}{\partial w} f(x_i, w)$$

$$= 0 \Rightarrow w^* = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

rare in ML. Usually a lot!  $x$  isn't usually a scalar but a vector of values.

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

$$\bar{x} = (x_1, x_2, \dots, x_d)^\top \Rightarrow d - \text{dimensional vector space}$$

$$y = f(x) = \sum_{j=1}^d w_j x_j + b$$

$$= x^\top \bar{w} + b$$

$$\tilde{x} = \begin{pmatrix} \tilde{x}_1 \\ \vdots \\ \tilde{x}_n \end{pmatrix}, \quad \tilde{w} = \begin{pmatrix} \bar{w} \\ b \end{pmatrix}$$

$$y_i = f(\tilde{x}) = \tilde{w}^\top \tilde{x}$$

$$e_i = y_i - \tilde{w}^\top \tilde{x};$$

$\downarrow$  error

objective

$$E(\tilde{w}) = \sum_{i=1}^N (y_i - \tilde{w}^\top \tilde{x})^2$$

$$\tilde{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \quad \tilde{X} = \begin{bmatrix} \tilde{x}_1^\top & \vdots & 1 \\ \tilde{x}_2^\top & \vdots & 1 \\ \vdots & \ddots & 1 \end{bmatrix} \quad N \text{ rows}$$

$d+1$  columns

$$\begin{aligned} E(\tilde{w}) &= \|\tilde{y} - \tilde{X}\tilde{w}\|_2^2 \quad \rightarrow \begin{cases} \text{two normed square} \\ \sqrt{\sum \text{of the elements}} \end{cases} \\ &= (\tilde{y} - \tilde{X}\tilde{w})^\top (\tilde{y} - \tilde{X}\tilde{w}) \\ &= \tilde{y}^\top \tilde{y} - 2\tilde{w} \quad \left| \begin{array}{l} \|v\|_2^2 = \sum v_j^2 = \tilde{v}^\top v \end{array} \right. \end{aligned}$$

Input:  $\vec{x} \in \mathbb{R}^d$ ,  $\vec{x} = [x]$

Note: Assign is out

Output:  $y \in \mathbb{R}$ ,

Eg. predict your grade using previous grades.

Start it already!

Eg. predict expected weight of child

Possible inputs  
 - weight  
 - height  
 - etc ..

$x$  is features, measurements of the world which will base your predictions

$y = f(\vec{x}) = \tilde{w}^T \vec{x}$  ← model fitting hyper plane to data.

If  $f$  is a vector instead of a scalar.

$$\vec{y} = f(\vec{x}) = \tilde{W}^T \vec{x}$$

$$\tilde{w} \in \mathbb{R}^{d+1 \times K}$$

$$\begin{aligned}\tilde{W} &= [\tilde{w}_1 \dots \tilde{w}_K] \\ &= \begin{bmatrix} \tilde{w}_1 & \tilde{w}_2 & \dots & \tilde{w}_K \\ b_1 & b_2 & \dots & b_K \end{bmatrix}\end{aligned}$$

Data:  $\{(\vec{x}_i, \vec{y}_i)\}_{i=1}^N$

$j^{\text{th}}$  element of  $i^{\text{th}}$  training sample

$$\hat{y}_{ij} = \tilde{w}_j^T \vec{x}_i \rightarrow y_{ij}$$

$$= \sum_{i=1}^N \sum_{j=1}^K (y_{ij} - \tilde{w}_j^T \vec{x}_i)^2$$

loss and error are interchangeable in this

the loss is non-negative. when loss = 0, you get a perfect solution.

$\hat{y}_j$  be  $j^{\text{th}}$  output for all  $N$  training points

per.

$y$  - describes colours  
 $y'$  - describes rows

$$\begin{bmatrix} \bar{y}_1 & \dots & \bar{y}_n \end{bmatrix}$$

$\tilde{x} = [\tilde{x}_1 \ \tilde{x}_2 \ \dots \ \tilde{x}_n]^T$  ← with transpose on  
rows matrix  $X$ .

$$E[\tilde{w}] = \sum_{j=1}^K \|g_j' - \tilde{x}\tilde{w}_j\|^2$$

↑ output vector matrix  
↑ weights

, shows  $K +$   
squared problems

$$E(\tilde{w}) = \|y - \tilde{x}\tilde{w}\|_F^2$$

Frobenius norm

$$y = [\bar{y}_1 \ \bar{y}_2 \ \dots \ \bar{y}_n]$$

Frobenius norm

normal vector func norm applied to matrix

$$\|A\|_F^2 = \sum_i \sum_j a_{ij}^2$$

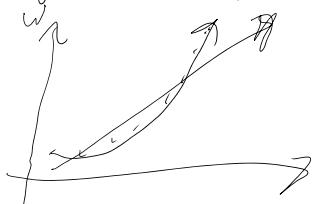
Two norms

$$\|A\|_2^2 = \sum_i \sum_j a_{ij}^2$$

## Non linear regression

when you get data plot it, find a way to visualize it.

height vs weight



Think relation between feature and outputs

non-linearity make life difficult. But class of function easy to solve: Basis Function regression

## Basis Function Regression

$$u = f(x) = \sum_{k=1}^K w_k b_k(x)$$

linear function regression

$$y = f(x) = \sum_{k=1}^K w_k b_k(x)$$

— use sin

— use poly nomial

— use linear

## Two basis functions

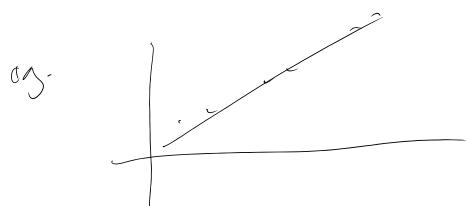
### ① polynomials

$b_k(x) = x^k \leftarrow x^k$  is monomial because it has only 1 order

Weighted sum of monomial is polynomial.

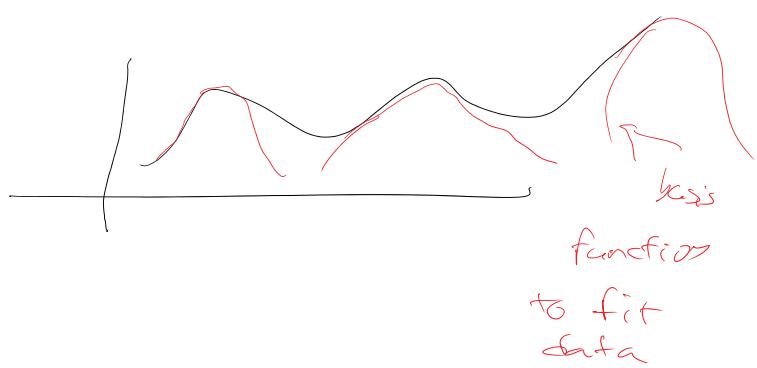
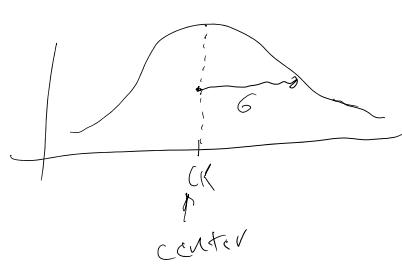
$$f(x) = w_0 + w_1 x + w_2 x^2 + \dots$$

1<sup>st</sup> assignment - fit polynomial function to data.



### Radial Basis Function : (RBF)

$$b_k(x) = e^{-\frac{(x-c_k)^2}{G_k^2}} \leftarrow \text{looks similar to normal}$$



$$y = f(x) = \sum w_k b_k(x)$$

$$\sum (x_i, y_i) \}_{i=1}^N$$

## Estimation

$$E(\bar{\omega}) = \sum_i (y_i - f(x_i))^2$$

$$\bar{\omega} = (\omega_1, \dots, \omega_K)^T \quad , \quad \bar{y} = (y_1, \dots, y_N)^T$$

$$B = [B_{ik}] = B_{ik} = b_{ik}(x_i)$$

$$= \begin{bmatrix} b_1(x_1) & \dots & b_K(x_1) \\ b_1(x_2) & b_2(x_2) & \dots & b_K(x_2) \\ \vdots & & & \\ b_1(x_n) & \dots & \dots & b_K(x_n) \end{bmatrix}$$

feature set for  $x$ .  
basis function vary  
each row

↑  
all the inputs vary

$$\begin{aligned} E(\bar{\omega}) &= \| \bar{y} - B\bar{\omega} \|^2 \\ \text{objective} &= (\bar{y} - B\bar{\omega})^T (\bar{y} - B\bar{\omega}) \\ &= \end{aligned}$$

how do we find  
optimal  $\bar{\omega}$ ?  
Take gradient of  $E$

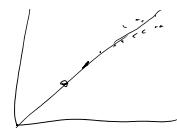
$$\nabla E = \frac{\partial E}{\partial \bar{\omega}} = 0$$

$$\Rightarrow \bar{\omega}^* = (B^T B)^{-1} B^T \bar{y}$$

↑  
(diag +  
scaled)  
inverse

## Reasons for overfitting

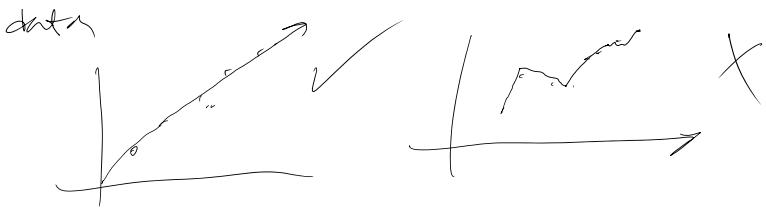
- Too many parameters and too few data.
- modelling the noise process if self.
- model uncertainty



## Regularization

Smooth models are better than non smooth





if you got a good fit using 4 basis function, it's better than 20 basis functions.

$$E(\bar{w}) = \|\bar{y} - B\bar{w}\|^2 \leftarrow \text{objective function}$$

$$= \underbrace{\|\bar{y} - B\bar{w}\|^2}_{\text{data term}} + \lambda \|\bar{w}\|^2 \leftarrow \begin{array}{l} \text{extra terms} \\ \text{to penalize } w \text{'s} \\ \text{that are big} \\ \left( \sum w^2 \right) \end{array}$$

regulation parameter  
 $L_2$  controls balance

smoothness term

$$y = f(x) = w_1 x + w_2 x^2 + w_3 \dots$$

$$\frac{\partial f}{\partial x} = w_1 + 2w_2 x + \dots$$

$$\frac{\partial f}{\partial x} = 0 + 2w_2 + \dots$$

$$\begin{aligned} E(\bar{w}) &= \|\bar{y} - B\bar{w}\|^2 + \lambda \|\bar{w}\|^2 \\ &= (\bar{y} - B\bar{w})^T (\bar{y} - B\bar{w}) + \lambda \bar{w}^T \bar{w} \\ &= \bar{w}^T B^T B \bar{w} + \lambda \bar{w}^T \bar{w} - 2 \bar{w}^T B^T \bar{y} + \bar{y}^T \bar{y} \\ &= \bar{w}^T (B^T B + \lambda I) \bar{w} - 2 \bar{w}^T B^T \bar{y} + \bar{y}^T \bar{y} \end{aligned}$$

$$\nabla E = 0$$

$$2(B^T B + \lambda I)\bar{w} - 2B^T \bar{y} = 0$$

$$\therefore \bar{w}^* = B^T B + (\lambda I)^{-1} B^T \bar{y}$$

### Ridge regression

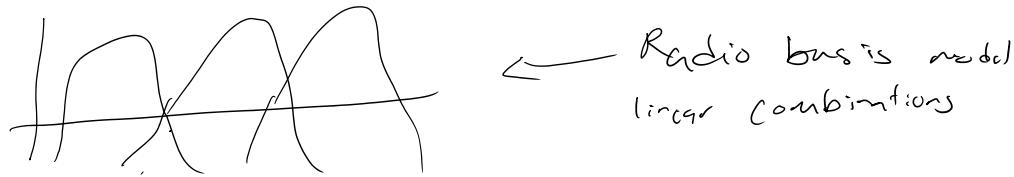
form of bias regression using regularization process to encourage your model to be smooth.

## Regression

We have been using parametric models

parametric vs non parametric model?

parametric models have fixed # of parameters. # of parameters in the model does not depend on the training set. You don't have to remember the training data.



radio basis model is taking data in the local neighbourhood and fitting it with a bump.

Non-parametric model ex/

$$\{(x_i, y_i)\}_{i=1}^N$$

### K-nearest neighbour Regression (k-NN Regression)

Given a test point  $x$  get the nearest points to the  $x$  and average them.

$$y = \frac{1}{k} \sum_{i \in N_k(x)} y_i$$

set of neighbours Indices into training data given point  $x$

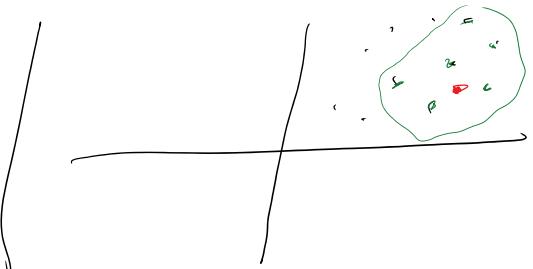
We don't usually choose  $k=1$  because data is noisy

Think of  $k$  as a parameter like spacing of rbf functions

$k \sim$  hyperparameter.

hyperparameter of polynomial - power/order

larger value of  $k$ , smoother the fit



- training set

- training point

• - nearest neighbor  
point, average  
from

## Estimation Theory, Bayes' optimal form of data fitting and prediction

### Uncertainty

Prob: assign beliefs to events without observing them.

A: die shows 3

$$P(A) = \frac{|A|}{|S|} = \frac{1}{6}$$

B: die shows 1

$$P(B) = \frac{1}{6}$$

C: dice sum to 8, 2 dice

$$P(C) = \frac{|C|}{|B|} = \frac{6}{6^2} = \frac{5}{36}$$

$\begin{matrix} 5, 3 \\ 3, 5 \\ 2, 6 \\ 6, 2 \\ 4, 4 \end{matrix} \Bigg) \quad |C|=5$

$$P(A \cap C) = P(A \cap C)$$
$$= \frac{1}{36}$$

$$P(B \cap C) = 0$$

$$P(C | A) = \frac{P(A \cap C)}{P(A)} = \frac{\frac{1}{36}}{\frac{1}{6}} = \frac{1}{6}$$

$$\perp P(A, B) = P(A) P(B)$$

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A) = \sum_i P(A, B_i) \leftarrow \text{Joint marginal.}$$

Hedroschmidt: some data points have more noise than others.



## Bayes Rule

$P(M | D)$  → A condition on the data  
 ↑  
 model      Data

"Posterior distribution" — probability distribution  
 of the model after  
 you seen the  
 data

$\Rightarrow$  allows you to say which models are conditioned to my data.

to my data.

likelihood: replicability you will see the data with the given data

Bayes rule says  $P(M|D) = \frac{P(D|M) p(M)}{P(D)}$

an prior - by flipping coin  
prior belief is  
 $\frac{1}{2}$

an evidence

Inference: In ML, computing prob dist'n unknown parameter  
interest. Estimating prob distribution over unknown

$p(m|D)$

Estimation: estimating single model from data  
"best model"

1) MAP (maximum a posteriori) Estimation (modo)

A model that maximizes

$$m = \omega \cdot (\text{color}(n) + \text{color}(n))$$

A model that maximizes

$$\Theta_{MAP} = \underset{\theta}{\text{argmax}} p(\theta | D) = \underset{\theta}{\text{argmax}} p(D|\theta)p(\theta)$$

parameter  
of model

specific instance of the model class which  
maximize posterior distribution

2) ML (Maximum likelihood) Estimation

$$\Theta_{ML} = \underset{\theta}{\text{argmax}} p(D|\theta)$$

Difference between map and likelihood is we  
are not using prior belief

Note: equal if uniform distribution.

---

N coin flips:  $c_{1:N} = (c_1, c_2, \dots, c_N)$

$c_1, \dots, c_n$  - outcome

flip 1: N times.

$$\Theta : p(H) \quad \text{equivalently} \quad \begin{aligned} p(c=H) &= \Theta \\ p(c=T) &= 1 - \Theta \end{aligned}$$

parameter

prior

$$\begin{aligned} p(\theta) &= 1 \\ \theta &\sim U(0, 1) \end{aligned}$$

assume  $\perp$

$$p(c_{1:N} | \theta) = \underbrace{\prod_{i=1}^N p(c_i | \theta)}$$

$$p(\theta | c_{1:N}) = \frac{p(c_{1:N} | \theta) p(\theta)}{p(c_{1:N})}$$

estimation:

$m^*$

$r_1 \sim 1 \sim r_2 \sim 1$

estimation:

$$\begin{aligned}\hat{\theta}^* &= \arg \max_{\theta} p(c_{1:N} | \theta) p(\theta) \\ &= \arg \max_{\theta} \log (p(c_{1:N} | \theta) p(\theta)) \quad \left. \begin{array}{l} \text{use these} \\ \text{interchangeable} \end{array} \right\} \\ &= \arg \min_{\theta} (-\log p(c_{1:N} | \theta) p(\theta))\end{aligned}$$

$$N = 1000$$

$$\# \text{ of heads} = 750$$

then

$$\# \text{ of tails} = 250, \text{ assume no coin landed on tails}$$

$$p(c_{1:N} | \theta) = \prod_{i=1}^{1000} p(c_i | \theta) = \theta^{750} (1-\theta)^{250}$$

$$c_{1:5} = \text{HTHTT}$$

$$\begin{aligned}p(c_{1:5} | \theta) &= p(c_1 | \theta) p(c_2 | \theta) p(c_3 | \theta) \dots \\ &= \theta \theta (1-\theta) (1-\theta) \\ &= \theta^3 (1-\theta)^2\end{aligned}$$

$$p(\theta | c_{1:1000}) = k \underbrace{p(c_{1:1000} | \theta)}_{\theta^{750} (1-\theta)^{250}} \underbrace{p(\theta)}_1$$

$$\simeq K \theta^{750} (1-\theta)^{250}$$

# minimize neg log posterior

$$= -\log(K) - 750 \log(\theta) - 250 \log(1-\theta) \rightarrow \text{proportion dawg}$$

Minimize it: # take derivative, set it = 0

$$\frac{d}{d\theta} (p(\theta | c_{1:100}))$$

$$0 - \frac{750}{\theta} + \frac{250}{1-\theta} = 0$$

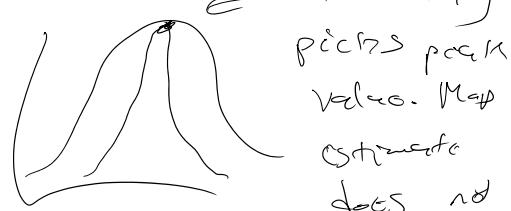
$$\frac{750}{\theta} = \frac{250}{1-\theta}$$

$$750(-\theta) = 250(\theta)$$

$$750 = 250(\theta) + 750(0)$$

$$\frac{750}{1000} = \theta^* \text{ or } \theta_{\text{map}} \text{ frequency of occurrence.}$$

This says map estimate =  $\frac{\bar{H}}{N}$



$$E[\theta] = \int \theta p(\theta | c_{1:100}) d\theta$$

$$= \int \theta K \theta^{750} ((-\theta)^{250} \times 1)^{100} d\theta$$

⋮

$$= \frac{750+1}{750+250+2}$$

$$= \frac{751}{1002} \quad \text{In general} = \frac{H+1}{N+2} \quad \begin{cases} \text{Bayes estimate} \\ (\text{mean}) \end{cases}$$

Regression

← d-dimensional weights: treat as r.v.

/  $u \in \mathbb{R}$

## Regression

$$y = \bar{w}^T \bar{x} + \varepsilon \quad \leftarrow \text{model}$$

d-dimensional weights: treat as r.v.

scalar

d-dimensional set of measurements

source of noise

$$\begin{cases} y \in \mathbb{R} \\ \bar{x} \in \mathbb{R}^d \\ \bar{w} \in \mathbb{R}^d \\ \varepsilon \sim N(0, \sigma^2) \\ \underbrace{\qquad\qquad}_{\text{uncertainty in}} \\ \text{measurements} \end{cases}$$

$$\text{mean of } y: \bar{w}^T \bar{x}$$

$$\text{variance: } \sigma^2$$

$y \sim$  variable that is gaussian

$$y | \bar{w}, \bar{x} \sim N(\bar{w}^T \bar{x}, \sigma^2)$$

$$p(y | \bar{w}, \bar{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y - \bar{w}^T \bar{x})^2}{2\sigma^2}}$$

$\downarrow$  dimension is  
0 because of different weights

$$p(\bar{w}) \quad \text{assume } \bar{w} \sim N(0, \alpha I)$$

$\downarrow$  mean variance

Isotropic prob distribution is the same direction

$$\begin{aligned} p(\bar{w}) &= \prod_{i=1}^d p(w_i) \\ &= \prod_{i=1}^d \frac{1}{\sqrt{2\pi\alpha}} e^{-\frac{1}{2\alpha} w_i^2} \\ &= \left( \frac{1}{\sqrt{2\pi\alpha}} \right)^d \frac{1}{c^{d/2}} e^{-\frac{1}{2\alpha} \bar{w}^T \bar{w}} \end{aligned}$$

Hand in written work in the drop box  
IC400 - Drawel CSC11

Hand in code online. Assignment due before lecture. Quiz during lecture.

Midterm: After the reading week

Date: October 20<sup>th</sup>, Saturday Afternoon.

No cheat sheet.

### Regression

$$y = \bar{x}^T \bar{w} + \epsilon$$

$$\bar{x} \in \mathbb{R}^d \quad \epsilon \sim N(0, \sigma^2 I)$$

$$\bar{w} \in \mathbb{R}^d$$

$$P(y | \bar{x}, \bar{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(y - \bar{x}^T \bar{w})^2}{\sigma^2}}$$

$$P(\bar{w}) \sim N(0, \frac{\alpha}{d} I)$$

covariance matrix which is isotropic with variance  $\alpha$

$$P(\bar{w}) = \frac{1}{(2\pi\alpha)^{\frac{d}{2}}} e^{-\frac{1}{2}\alpha \bar{w}^T \bar{w}}$$

### Posterior

$$D = \sum_{i=1}^N (\bar{x}_i, y_i)$$

$$P(\bar{w} | D) = \frac{P(y_{1:N} | \bar{x}_{1:N}, \bar{w}) P(\bar{w})}{P(y_{1:N})} \text{, assume } \bar{w} \perp \text{ of } y_{1:N}$$

assume noise is  $\perp$

$$\text{Independence } y_i = \bar{x}_i^T \bar{w} + \epsilon_i \quad \sum_i \epsilon_i = 0$$

$$P(\epsilon_1, \epsilon_2) = P(\epsilon_1) P(\epsilon_2)$$

$$\Rightarrow P(y_1, y_2 | \bar{x}_1, \bar{x}_2, \bar{w}) = P(y_1 | \bar{x}_1, \bar{w}) P(y_2 | \bar{x}_2, \bar{w})$$

$$P(\bar{w} | D) = \frac{\prod_{i=1}^N P(y_i | \bar{x}_i, \bar{w}) P(\bar{w})}{C}$$

$e$  raised to quadratic multiplied by  $e$  raised to the quadratic is sum of quadratic

$$P(y_{1:N} | \bar{x}_{1:N}, \bar{w}) = \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{x}_i^T \bar{w})^2}$$

$$P(\bar{w} | D) = C' \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{x}_i^T \bar{w})^2} \left( \frac{1}{2\pi\alpha} \right)^{\frac{d}{2}} e^{-\frac{1}{2\alpha} \bar{w}^T \bar{w}}, \text{ take } -\log$$

$\log$  of product is sum of  $\log$

$$-\log P(\bar{w} | D)$$

$$= -\frac{1}{2} \log C + \frac{N}{2} \log \left( \frac{1}{2\pi\sigma^2} \right) + \frac{d}{2} \log \left( \frac{1}{2\pi\alpha} \right) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{x}_i^T \bar{w})^2 + \frac{1}{2\alpha} \bar{w}^T \bar{w}, \text{ disappears}$$

$$\begin{aligned}
 & -\log p(w|x) \\
 & = -\cancel{\log \sigma} + \frac{N}{2} \cancel{\log(2\pi\sigma^2)} + \frac{1}{2} \cancel{\log(2\alpha)} + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{x}_i^T \bar{w})^2 + \frac{1}{2\alpha} \bar{w}^T \bar{w}, \quad \text{disappear because it doesn't depend on } w. \\
 & -\log p(\bar{w}|x) = k + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{x}_i^T \bar{w})^2 + \frac{1}{2\alpha} \bar{w}^T \bar{w}
 \end{aligned}$$

#Take derivative and set it to zero

#multiply by a constant

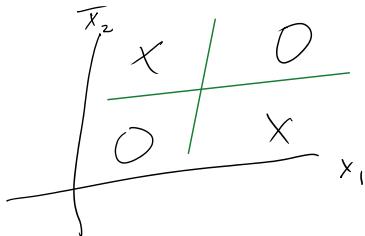
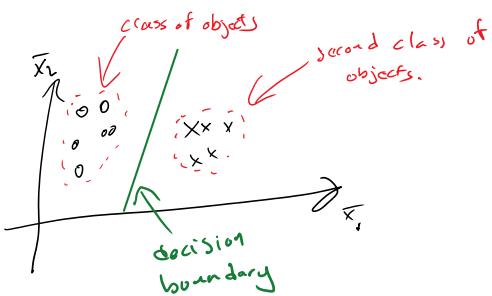
$$\begin{aligned}
 & = k' + \sum_{i=1}^N (y_i - \bar{x}_i^T \bar{w})^2 + \left( \frac{2\sigma^2}{2\alpha} \right) \bar{w}^T \bar{w} \\
 & \quad \text{, form of ridge regression} \\
 & \lambda = \frac{2\sigma^2}{2\alpha}
 \end{aligned}$$

What is a gaussian?

What is a prior?

Classification

$$\begin{aligned}
 y &= f(x) \\
 y &\in \{-1, 1\} \\
 y &\in \{1, \dots, k\}
 \end{aligned}$$



How to determine if  $X$  or  $O$ ? Guess based on boundaries

Complexities are based on the model.

linear separable: A linear function can separate the data.

Classification by regression:

$$\text{data: } \{(x_i, y_i)\}_{i=1}^N, y_i \in \{-1, 1\}$$

#Find the best separating hyperplane

$$\underset{w}{\text{weights}} \rightarrow w^* = \arg \min_{\bar{w}} \sum_{i=1}^N (y_i - \bar{x}_i^T \bar{w})^2 \quad \leftarrow \text{squared error}$$

$$f(\bar{x}) = \bar{w}^T \bar{x}$$

The classifier is the sign of  $f(x)$

$$\text{classifier } \text{sgn}(f(x)) = \text{sgn}(\bar{w}^T \bar{x})$$

$$w^T(x) > 0, \text{ above hyper plane}$$

$\omega^T(x) < 0$ , below the hyper plane

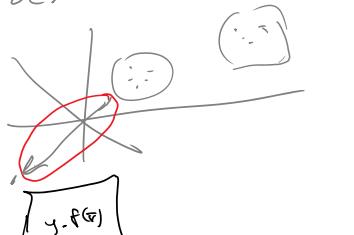
$$\text{sign}(z) = \begin{cases} -1 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

This works well in only certain situations

- points closely clustered together
- spread is more or less the same



Doesn't work well for



since loss function  
is huge



$$\frac{0-1 \text{ loss}}{L_{0-1}(y_i, f(x_i))} = \frac{\text{if } y_i \neq f(x_i)}{1}$$

$\sum_i L_{0-1}(y_i, f(x_i)) \leftarrow$  loss function that  
counts # of mistakes

problem? hard to optimize

Suppose  $f(x) = -2$

$$y = 1$$

$$f(x) \times (y) = -2 \leftarrow \text{error.}$$

penalize you for confident and right and  
confident & wrong.

### K-NN Classifier

Find the K nearest neighbours to test inputs (with indices in  $N_k(x)$ )

Let  $y = \text{sign}\left(\sum_{i \in N_k(x)} y_i\right) \leftarrow$  This is a way of voting. Sum votes and look at the sign.

$$w(v_i) = e^{-\frac{1}{2\sigma^2} \|x - x_i\|^2}$$

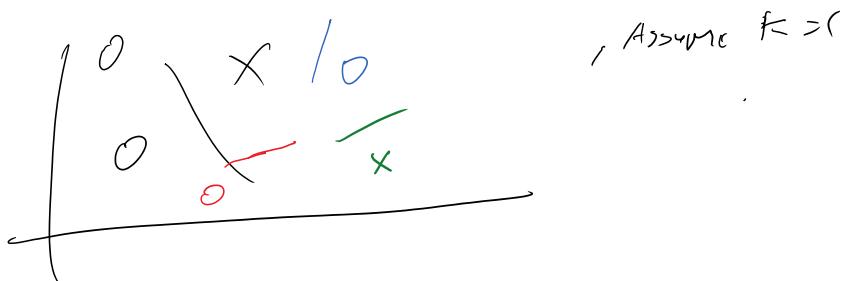
# Question on midterm how do  
K-NN classifier well.

Bad: search, lots of point then you have to search all the points.

This is non parametric method, no fixed parameter. Algorithm depends on the amount of data.  
parametric model: # of parameter doesn't change.

non parametric would do feature learning.

what does decision boundary look like



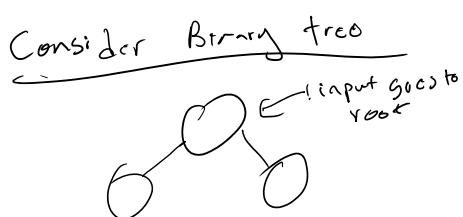
Decision boundary can be complex and can overfit.

larger value of  $k \leftarrow$  # of nearest neighbour  
 $\Rightarrow$  blurring out decision boundary

Decision Tree  $\hookrightarrow$  go to, often used in practice

stump  $\rightarrow$  tree  $\rightarrow$  forest

Decision tree can be used for regression or classification.



$$\begin{cases} M \text{ internal nodes} - \text{split nodes} \\ m+1 \text{ leaf nodes} \\ \text{at internal node } j \\ t_j(x) : \mathbb{R} \rightarrow \{-1, 1\} \end{cases}$$
$$T_j(x) = \begin{cases} -1, & \text{left branch} \\ 0, & \text{take right branch} \end{cases}$$

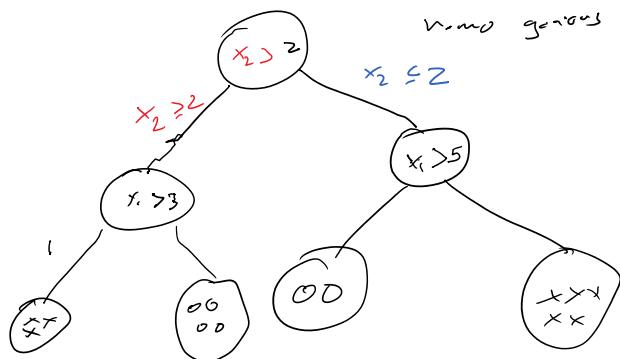
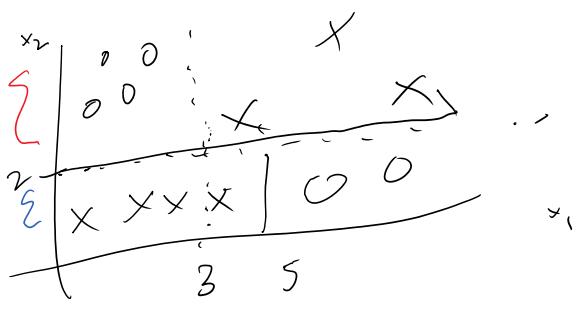
1. input goes to root of tree

2. gonna ask a feature eg  $x \geq 2$  or  $2 \leq x$

then go to the new node.

3. New  $t > , t \geq 3$  or  $t \leq 3$  then go down the new node.

goal of split: Separate class of data. Split it for each it more generic



Once you get to leaf node, take a vote.

# of regions grows exponentially depending on depth of the tree.

⑩ leaves of tree define

$$p(c|\text{node } j) = \frac{N_{jF}}{N_j} \quad \left. \begin{array}{l} \text{fraction of data points that reach} \\ \text{that node that have class } c. \end{array} \right.$$

$N_j$  = # points reach node  $j$

$N_{jF}$  # points in node  $j$  with class  $c$

Midterm covers everything up to today.

October 20<sup>th</sup> 5:00 - 6:00.

### Decision Tree

Sequence of decisions, leaf nodes will give you prob distribution over classes.

Internal nodes  $t_j(\bar{x}) : \mathbb{R}^d \rightarrow \{-1, 1\}$

leaf nodes,  $\Leftrightarrow$  node  $j$  with

$N_j$  items    prob( $y=c | \text{node } j$ )

leaf node  $j$ , with  $N_j$  items

$N_{j,c}$  with class  $C$

$$P(y=c | j) = \frac{N_{j,c}}{N_j}$$

### Decision Trees

node  $j$      $t_j(\bar{x})?$   
split function

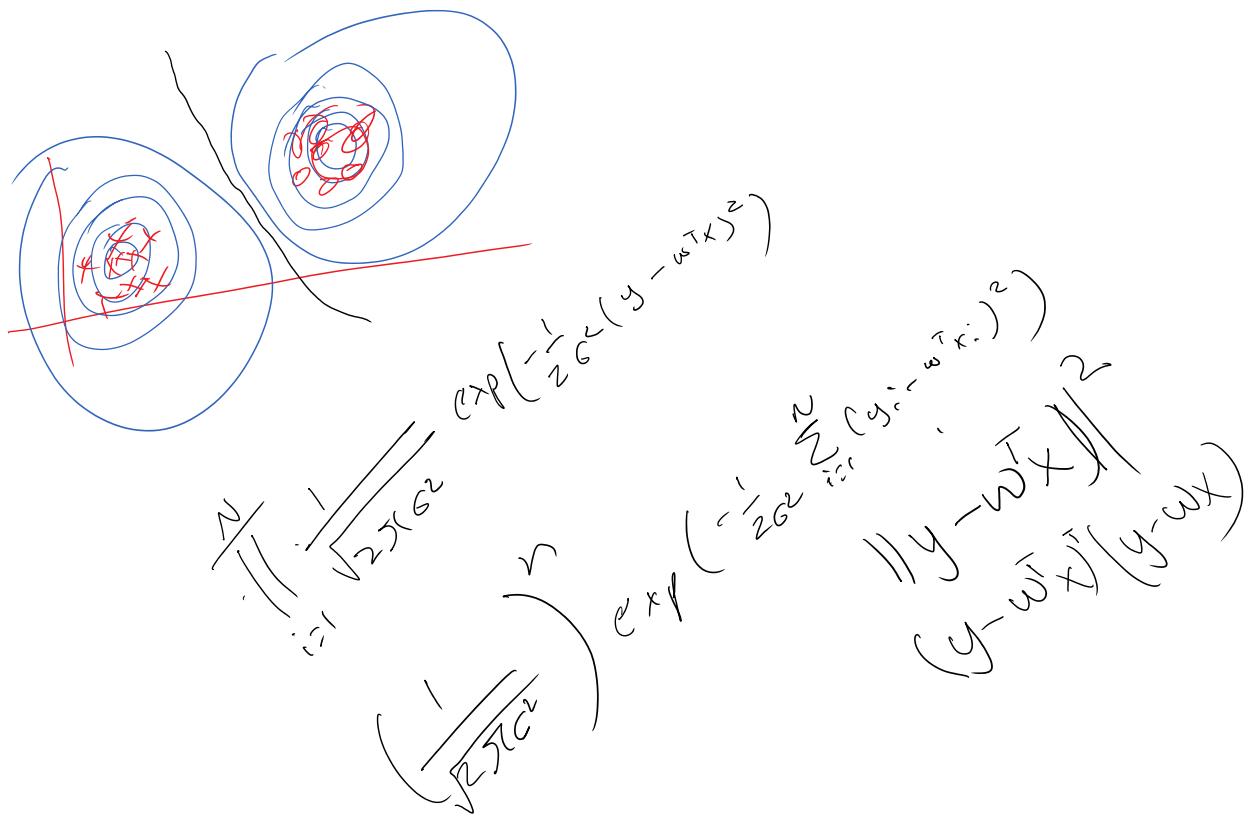
Given  $m$ -items  $\{\bar{x}_i\} = D_j$     data at node  $j$   
(to)



what will the split function look like?

Assume univariate split

↳ split function that only depends on one variable



Only 3 assignments!

Assignment is getting extended!

↳ Monday or Tuesday

Midterm marks within next few days.

Next week in tutorial go over questions &amp; solutions.

K nearest neighbor: can store all data  
- classifying all data

Decision Trees: many competitions in ML are won by random forest classifiers

Random Forest: any of decision trees

Class Conditional model:

Naive Bayes: assumption features on class is ⊥

Logistic Regression:

5 types of classifier

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

↓  
solving this  
using → discriminative  
model

↑  
generative model  
If you know the class  
you can generate data  
from it.

Discriminative model: give me measurement, give you back classes. Eg Decision tree of height, width, colour, texture etc..

10 measurements and run decision tree algorithm. Assume works with 2 measurements, therefore the other measurements, therefore the other 8 measurements is unused

Generative model does more work and solving a harder problem than you need to

Discriminative models are usually better.

Logistic regression, assume  $p(y=1|x) = \frac{1}{1+e^{-w^T x}}$   
presuppose parametric forms.

Gaussian class condition - generative model

Naive Bayes - generative model.

How to choose a good model?

Simple way: take training set → break it up  
into training set and validation set.

Pick a validation set because suppose training data is like test data. Take 80% and 20%  
20% → a validation set and train only on 80%

How to choose validation set? At random

Sample without replacement from random.

Problem with validation set techniques

Lack of data. 1000 training data  
- 200 validation

1000



problem with validation set techniques

$$\begin{array}{c} \text{lack of data. } 1000 \text{ training data} \\ - 200 \text{ validation} \\ \hline 800 \text{ training left} \end{array}$$

less to train on. Weak when you don't have enough data.

when lack of data:

Validation: Train and choose split 80,20  
60,40 etc... You need enough for validation

### N-fold Cross Validation

Take data, rather than make 20 out for validation

do 5 random samples of 20. Breaking

do 5 random samples of 20. Breaking

partition of 5 subsets. Build 5 models

and test it on 20%. Average over 5 samples.

How big is N?  $N \uparrow$  when data set  $\downarrow$

limit  $\rightarrow$  LOOCV

$\rightarrow$  leave one out cross validation

Avg errors and that will tell you how well your model is

Let  $\hat{y}_i$  be the prediction of model trained on  $N-1$  pts (i.e. w/o  $i^{th}$ )

$$\text{LOOCV} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

### LS regression

$$w^* = (X^T X)^{-1} X^T \bar{y}$$

(compute optimal weights)

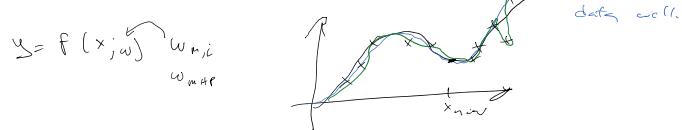
$$\hat{y} = Xw^* = \underbrace{X(X^T X)^{-1} X^T}_{H} \bar{y} = H\bar{y}$$

some matrix  $H$  times  $y$ .

$$\begin{aligned} \text{LOOCV} &= \frac{1}{N} \sum_i \text{MSE}_i \\ &= \frac{1}{N} \sum_i \left( \frac{y_i - \hat{y}_i}{h_{ii}} \right)^2, \quad h_{ii} \text{ is } i^{th} \text{ diagonal element} \\ &\quad \text{of } H, \text{ leverage of } i^{th} \text{ point.} \end{aligned}$$

### Bayesian Method

predict  $y_r$  given  $D = \{(x_i, y_i)\}_{i=1}^N$



$$p(w|D) = \frac{p(D|w)p(w)}{P(D)}$$

\* Find entire distribution  $p(w|D)$

\* find  $p(y_{\text{new}} | x_{\text{new}}, D)$

We never really care about  $w$  but more  $y$ .

$$p(y|D) = \int p(y, w|D) dw \quad \rightarrow y = f(x, w) \text{ is close to data}$$

$$\begin{aligned} p(y|x,D) &= \int p(y, w|x, D) dw \\ &= \int p(y|w, x, D) p(w|x, D) dw \quad \rightarrow p(\tilde{w}|D) \text{ far from data} \end{aligned}$$

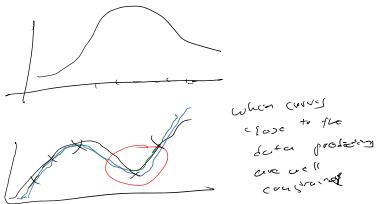


$$p(y|x) = \int p(y, \bar{w}|x, D) d\omega$$

$$= \int p(y|\bar{w}, x, D) p(\bar{w}|D) d\omega \quad \begin{array}{l} \text{uncertainty} \\ \downarrow \\ \text{uncertainty} \end{array}$$

2)  $p(\bar{w}|D)$  far from data

what you get for Bayesian  
is a whole probability  
distribution over  $y$ .



### Bayesian Regression / Gaussian process

$$y = \bar{w}^\top b(x) + n \leftarrow \text{random noise} \quad y \in \mathbb{R}$$

$$n \sim N(0, \sigma^2) \quad x \in \mathbb{R}$$

$$\bar{w} \sim \mathcal{N}(0, \frac{1}{\alpha} I_d) \quad b(x) \in \mathbb{R}^d$$

$$p(y_{1:N} | x_{1:N}, \bar{w}) = \prod_{i=1}^N p(y_i | x_i, \bar{w}) \quad \leftarrow \text{statistical independence}$$

Assume a prior,  $\bar{w} \sim N(0, \frac{1}{\alpha} I_d)$

$p(\bar{w}|D) = \frac{\text{likelihood prior}}{\text{posterior}}$

$= \frac{\prod_{i=1}^N p(y_i | x_i, \bar{w}) p(\bar{w})}{p(D)}$

$$-\log p(\bar{w}|D) = -\sum_{i=1}^N \log p(y_i | x_i, \bar{w})$$

$$= -\log p(\bar{w}) + \log p(D)$$

$$y_i | x_i, \bar{w} \sim N(b(x_i)^\top \bar{w}, \sigma^2) = \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^N}_{\text{variance}} (y_i - \underbrace{\bar{w}^\top b(x_i)}_{\text{mean}})^2 + \frac{\alpha}{2} (\bar{w}^\top \bar{w})^2 + c'$$

$$= \frac{1}{2\sigma^2} (\bar{y} - B\bar{w})^\top (\bar{y} - B\bar{w}) + \frac{\alpha}{2} (\bar{w}^\top \bar{w})^2 + c'$$

$$\bar{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \quad B = [b(x_1) \ b(x_2) \ \dots \ b(x_N)]^\top$$

$$\begin{aligned} &= \frac{1}{2\sigma^2} (\bar{y} - B\bar{w})^\top (\bar{y} - B\bar{w}) + \frac{\alpha}{2} (\bar{w}^\top \bar{w})^2 + c' \\ &= \frac{1}{2\sigma^2} (\bar{y}^\top - \bar{w}^\top B^\top)(\bar{y} - B\bar{w}) + \frac{\alpha}{2} (\bar{w}^\top \bar{w})^2 + c' \\ &= \frac{1}{2\sigma^2} (y^\top - \bar{y}^\top B^\top - \bar{w}^\top B^\top y + \bar{w}^\top B^\top B\bar{w}) + \frac{\alpha}{2} \bar{w}^\top \bar{w} + c' \\ &= \frac{1}{2\sigma^2} (y^\top y - \bar{y}^\top B^\top - \bar{w}^\top B^\top y + \bar{w}^\top B^\top B\bar{w}) + \frac{\alpha}{2} \bar{w}^\top \bar{w} + c' \end{aligned}$$

control drugs with

$$\begin{aligned}
 & \stackrel{26'}{=} \frac{1}{2\sigma^2} \left( \bar{y}^\top \bar{y} - 2\bar{y}^\top B\bar{\omega} + \underbrace{\bar{\omega}^\top B^\top B \bar{\omega}}_{\text{complicated square}} \right) + \frac{\alpha}{2} \underbrace{\bar{\omega}^\top \bar{\omega} + c''}_{\text{constant}}
 \end{aligned}$$

$$\begin{aligned}
 & = \frac{1}{2} \left( \bar{\omega}^\top \left( \frac{B^\top B}{\sigma^2} + \alpha I \right) \bar{\omega} \right) + c'' \quad , \text{ complicate square} \\
 & = \frac{1}{2} \bar{\omega}^\top \left( \frac{B^\top B}{\sigma^2} + \alpha I \right) \bar{\omega} - \bar{y}^\top \frac{B\bar{\omega}}{\sigma^2} + c'' \\
 -\log p(\bar{\omega} | D) & = \frac{1}{2} \underbrace{(\bar{\omega} - \hat{\omega})^\top K^{-1} (\bar{\omega} - \hat{\omega})}_{\text{since now see this form}} \quad , \text{ where } \hat{\omega} = \underbrace{k B^\top \bar{y}}_{\sigma^2} \\
 & \quad \text{log prob quadratic} \\
 & \Rightarrow \bar{\omega} \sim \text{gaussian}
 \end{aligned}$$

$$\begin{aligned}
 & \text{where } \hat{\omega} \text{ is mean} \\
 & K \text{ is covariance} \\
 p(y_{\text{new}} | x_{\text{new}}, D) & = \int p(y_{\text{new}} | x_{\text{new}}, \bar{\omega}) p(\bar{\omega} | D) \\
 & \sim \underbrace{N(\bar{b}(x_{\text{new}})^\top, \sigma^2)}_{\text{N}(\bar{b}, K)} \quad \underbrace{N(\bar{\omega}, K)}_{\text{N}(\hat{\omega}, K)}
 \end{aligned}$$

$$y_{\text{new}} | x_{\text{new}}, D \sim N(\hat{w}^\top \bar{b}(x_{\text{new}}), \sigma^2 \underbrace{\bar{b}^\top K \bar{b}}_{\text{b}(x_{\text{new}})})$$

collect terms with  
w, put the rest  
in the constant

the

$\sim$  mean

$\int \sim$  covariance