

## Lecture 2 - KNN

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

To classify new observation

1. Define distance measure (Euclidean)
2. Find Nearest neighbor for all new data points
3. Label them w/ label of their  $k$  nearest neighbor (pick the majority)

$$p(y=c|x, k) = \frac{1}{k} \sum_{\substack{i \in N_k(x) \\ \sim}} I(y_i=c) \quad \text{where } I(e) = \begin{cases} 1 & \text{if } e \text{ is true} \\ 0 & \text{if } e \text{ is false} \end{cases}$$

$\underset{\sim}{\text{K-NN of vector } x}$

### K-NN w/ weighted Classification

Pick the weighted majority label. Weight is inversely proportional to distance (closer = bigger weight)

$$p(y=c|x, k) = \frac{1}{\sum_{i \in N_k(x)}} \sum_{i \in N_k(x)} \frac{1}{d(x, x_i)} I(y_i=c) \quad d(x, x_i) = \text{distance b/w } x \text{ and } x_i$$

$$\hat{y} = \arg \max_c p(y=c|x, k) \quad \hat{y} = \frac{1}{\sum_{i \in N_k(x)}} \sum_{i \in N_k(x)} \frac{1}{d(x, x_i)} y_i \quad \text{for regression}$$

•  $x_{\text{new}}$  classify as class •  
 b/c the 2 point  
 is much nearer  
 = more weight

Hyperparameter -  $k$  - pick  $k$  that generalise the best

$p(y=b|x) = \text{high}$   
 $p(y=g|x) = \text{not so high}$

Distance Measure -  $L_2$  norm (Euclidean) =  $\|u-v\|_2 = \sqrt{\sum (u_i-v_i)^2}$  |  $L_1$ -norm =  $\sum |u_i-v_i|$

$L_\infty$  norm =  $\max_i |u_i-v_i|$

Angle =  $\cos \theta = \frac{u^T v}{\|u\| \|v\|}$

Mahalanobis :  $\sqrt{(u-v)^T S^{-1} (u-v)}$   
 $S$  PSD and symmetric

### Measure Classification Performance

- Confusion Table

		Actual	
		T	F
Predicted	T	TP	FP
	F	FN	TN

↓ Fraud

• Accuracy =  $\frac{TP + TN}{TP + TN + FP + FN}$

• Precision =  $\frac{TP}{TP + FP}$

• Recall / Sensitivity =  $\frac{TP}{TP + FN}$

• Specificity =  $\frac{TN}{FP + TN}$

• FNR =  $\frac{FN}{FN + TP}$   
late

• FPR =  $\frac{FP}{FP + TN}$

F1-Score =  $\frac{2 \cdot \text{precision} \cdot \text{sensitivity}}{\text{precision} + \text{sensitivity}}$

## Scaling issue

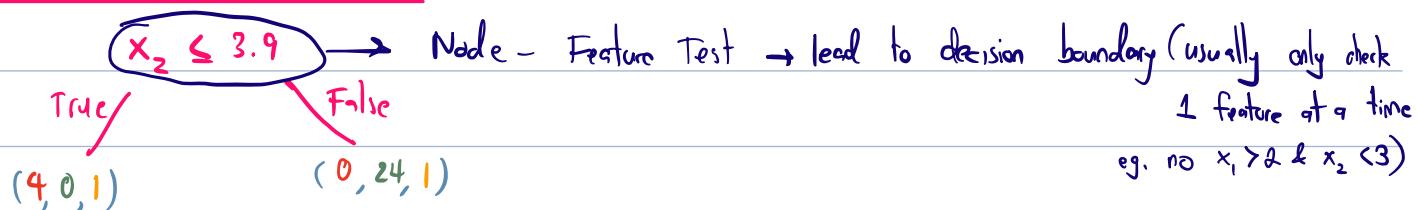
Data standardization - Scale each feature (col) to zero mean and unit variance

$$x_i = \frac{x_i - \mu_i}{\sigma_i}$$

## The Curse of Dimensionality

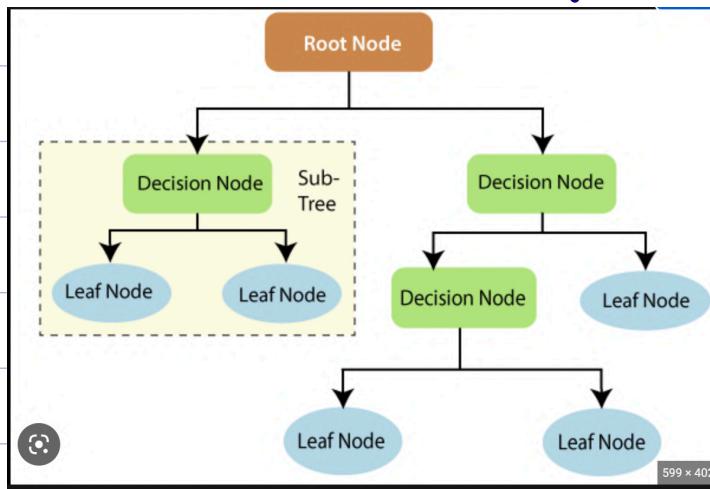
The higher the dimension, the further away the NN are. If the NN are too far away, then there's no notion of neighbor and NN doesn't work well.  
# of data has to grow exponentially with the # of feature (dim)

## Lecture 2: Decision Tree



Branch - Subsection of decision tree

Leaf - Node that doesn't split any further. Produce prediction or distribution



## To classify new sample $x$

1. Find region  $R$  that contain  $x$  and get the class distribution  $n_R = (n_{c_1,R}, n_{c_2,R}, \dots, n_{c_k,R})$

2. The prob that a data point  $x \in R$  should be classified to class  $c$  in region  $R$  is

$$p(y=c|R) = \frac{n_{c,R}}{\sum_{c \in C} n_{c,R}}$$

3. Sample  $x$  class = class that is the most common in the corresponding region

$$\hat{y} = \arg \max_c p(y=c|x) = \arg \max_c p(y=c|R) = \arg \max_c n_{c,R}$$

If region  $= (r, g, b) = \{4, 10, 6\} \rightarrow p(r) = \frac{4}{4+10+6} = \frac{4}{20} = 0.2$

$$p(b) = \frac{10}{20} = 0.5 \quad p(g) = \frac{6}{20} = 0.3$$

Naive Idea - Build all possible tree and test it on validation dataset to evaluate performance

### Heuristic

→ prob of picking class  $c$   
in dataset  $t$

Impurity measure let  $\pi_c = p(y=c|t)$   $i(t) = 0 \rightarrow$  node is pure  $i(t) = 1 \rightarrow$  class are equally distributed

1. Misclassification Rate  $\rightarrow i_E(t) = 1 - \max_c \pi_c$

2. Entropy  $\rightarrow i_H(t) = - \sum_{c \in C} \pi_c \log_2 \pi_c$  (let  $\log_2 0 = 0$ )

3. Gini Index  $\rightarrow i_G(t) = \sum_{c \in C} \pi_c (1 - \pi_c) = 1 - \sum_{c \in C} \pi_c^2$  Measure how often a randomly chosen instance would be misclassified if it was classified according to class distribution

The improvement when performing split feature  $s$  of dataset  $t$  into  $t_R$  and  $t_L$  for  $i(t) = i_E(t)$

$$\Delta i(s, t) = i(t) - p_L \cdot i(t_L) - p_R \cdot i(t_R)$$

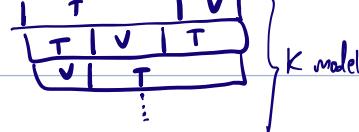
↑  
relative size of left & right region compare to before split

$$p_L = p_R = \frac{\text{data after latest split}}{\text{total data before the latest split}}$$

### Overfitting

- Poor generalisation
- Training performance increase w/ every split, Validation error got higher.
- Validation perf. tells us how well our model generalise, not training perf.

## K-fold Cross Validation



- Split learning data into K folds
- take one group as validation and use  $K-1$  folds for training
- Train model K times w/ diff division of training & validation data
- Avg the error over all model
- Try w/ diff setting of hyperparameter & model  $\rightarrow$
- Use all your training model and best hyperparameter for final training and testing of model

## LOOCV - Leave one out Cross Val..

- Train all but one sample  $\rightarrow$  if N sample  $\rightarrow$  N-fold cross validation  $\rightarrow$  train model N time

## Stopping Criteria (Pre-pruning)

- Distribution is pure  $i(t)=0$
- accuracy on validation set
- Max depth reach
- # samples in each branch below threshold  $t_n$
- benefit of splitting is below certain threshold  $\Delta i(s,t) < t_\Delta$

## Post Pruning $\rightarrow$ Reduce error pruning

Let  $T$  be our decision tree and  $t$  one of its inner nodes

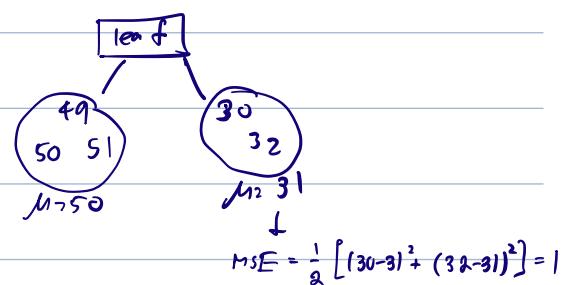
Pruning  $T$  wrt. to  $t$ :  $(T \setminus T_t)$  mean deleting all descendant nodes of  $t$  (but not node  $t$  itself)

- $\text{err}_{D_V}(T)$  = validation error of whole tree
- $\text{err}_{D_V}(T \setminus T_t)$  = validation error of tree w/o descendant of  $t$
- Prune tree at the node that reduce the validation error the most  $\rightarrow$  reduce overfitting

until for all nodes  $t$   $\text{err}_{D_V}(T) < \text{err}_{D_V}(T \setminus T_t)$

## DT for regression (if $y_i$ is real value rather than class)

- Compute mean at the leave
- Use MSE as splitting heuristic
- purity - all value are closed to the mean



## Ensemble

- Idea - Avg prediction of many (diverse) classifier to improve perf by reducing variance of the model by averaging
- Bagging - Create new dataset by sampling training set
- Train separate classifiers for each dataset
  - Combine predictions (Majority vote or average)

- Boosting
- Incrementally train (weak) classifier that correct previous mistakes
  - Give higher weight (Focus more) on misclassified example

- Stacking - Train meta-classifier w/ base classifiers' predictions as feature

## Bucket of nodes

### Random Forest • Bagging + Tree

- can be highly correlated
- reduce benefit of bagging since we need diverse classifier (for bagging)

- Idea - Use only a subset of randomly sampled data and feature to learn each tree  
(bagging at instance level (subset of sample) + bagging at feature level (subset of feature))

### Boosting : Ada Boost and XG Boost

Incrementally train weak learners (e.g. one-level-trees)

Initialise weight vector w/ uniform weight

Loop  $\rightarrow$  train weak learner on weighted example  $\rightarrow$  Increase weight for misclassified example

Predict the (error-based weighted) majority

## Lecture 3: Probabilistic Inference

Assume data is i.i.d (independent identical distribution)

Data and parameter  $\theta$  (could be probability of a coin toss)

Maximum Likelihood Estimator (MLE)

$$\hat{\theta}_{MLE} = \max_{\theta} p(D|\theta) \rightarrow \max_{\theta} \log(p(D|\theta)) \rightarrow \text{Point estimate}$$

find  $\theta$  that maximize the likelihood function  $p(D|\theta)$  (most likely to see observed data  $D$ )

The likelihood func is not prob distribution over  $\theta$  since  $\int p(D|\theta) d\theta \neq 1$  in general

MLE for any coin toss

$$\hat{\theta}_{MLE} = \frac{|T|}{|T| + |H|} \quad \text{find prob that next coin flip is } T \rightarrow F_{11} \sim B(\hat{\theta}_{MLE})$$

For Bernoulli, probability = parameter

$$F_{11} \sim B(\hat{\theta}_{MLE}) = p(F_{11} = T | \hat{\theta}_{MLE}) = \text{Ber}(F_{11} = T | \hat{\theta}_{MLE}) = \hat{\theta}_{MLE} = \frac{|T|}{|T| + |H|} = \frac{3}{3+7} = 30\%$$

But MLE doesn't consider prior knowledge  $\rightarrow$  Bayesian Inference

Bayesian Inference

$p(\theta)$  prior distribution represent our belief before we observe any data

Constraint for  $p(\theta)$

- $p(\theta)$  must not depend on the data right now
- $p(\theta) > 0$  for all  $\theta$
- $\int p(\theta) d\theta = 1$

$p(\theta|D)$  is the posterior distribution. Update our belief in the value of  $\theta$  after observing data

$$p(\theta|D) = \frac{p(D|\theta) p(\theta)}{p(D)}$$

$p(D|\theta)$  = likelihood func

$p(\theta)$  = prior distribution

$p(D)$  = evidence  $\rightarrow$  act like normalizing constant (so that  $\int p(\theta|D) = 1$ )

Posterior  $\propto$  Prior  $\cdot$  Likelihood

we can obtain evidence using sum rule of prob.  
 $p(D) = \int p(D,\theta) d\theta = \int p(D|\theta) p(\theta) d\theta$

Observing more data  $\rightarrow$  less effect on prior  $\rightarrow$  distribution lead toward current measurement more (closer to MLE)

If  $p(\theta) = 0 \rightarrow$  can't update posterior prob.  $\rightarrow$  always be 0

## Maximum a posterior estimation (MAP) & estimating Posterior Distribution (Fully Bayesian)

$$\text{MLE} - \theta_{\text{MLE}} = \arg \max_{\theta} p(D|\theta) \rightarrow \max_{\theta} \log p(D|\theta)$$

$\rightarrow$  return  $\theta$  that give the highest likelihood to see data D

$\rightarrow$  ignore prior belief & perform poorly if little data available

$$\text{MAP} - \theta_{\text{MAP}} = \arg \max_{\theta} p(\theta|D) \rightarrow \max_{\theta} \log p(\theta|D)$$

in the distri.

$\rightarrow$  return  $\theta$  - the maximizer of the posterior distribution (most probable  $\theta$  given the data D)

$\rightarrow \theta_{\text{MAP}}$  is a point estimate and  $p(\theta_{\text{MAP}}|D)$  is a highest probability (mode) in posterior distribution

## Fully Bayesian Approach - Compute the posterior distribution instead of just mode $\theta_{\text{MAP}}$

With an entire distribution, we can answer question such as

- Where is the mean under  $p(\theta|D)$   $\rightarrow$  not at maximizer b/c not symmetric

- Variance of  $p(\theta|D)$

Goal: get  $p(\theta|D)$  dist. by finding suitable  $p(\theta)$  and normalizing constant  $p(D)$

either by brute force or pattern matching

$$\int_0^1 p(\theta|D) d\theta$$

conjugate prior

choose  $p(\theta)$  which is a conjugate prior of the likelihood function  $\rightarrow$  The resulted Posterior Dist.  
(Pattern Matching)

will be in the same family as the prior  
(might not be the same distri as prior always)

e.g. likelihood =  $p(D|\theta) \sim \text{Ber}(\theta) = \theta^{\lvert T \rvert} (1-\theta)^{1+\lvert T \rvert}$

Conjugate prior distribution =  $p(\theta) = p(\theta|a,b) \sim \text{Beta}(a,b)$

So  $P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)} \propto p(D|\theta) \cdot p(\theta) = \theta^{\lvert T \rvert} (1-\theta)^{1+\lvert T \rvert} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$

$$\propto \theta^{\lvert T \rvert + a - 1} (1-\theta)^{1 + \lvert T \rvert + b - 1}$$

$$\theta_{\text{MAP}} = \arg \max_{\theta} p(\theta|D) = \arg \max_{\theta} \theta^{\lvert T \rvert + a - 1} (1-\theta)^{1 + \lvert T \rvert + b - 1} = \frac{\lvert T \rvert + a - 1}{1 + \lvert T \rvert + a + b - 2}$$

$$p(\theta|D) = \theta^{\lvert T \rvert + a - 1} (1-\theta)^{1 + \lvert T \rvert + b - 1}$$

$\rightarrow$  not normalise

find  $p(D) \Rightarrow$  pattern matching

$$\text{Beta}(\theta | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \rightarrow \text{already normalise}$$

So we know that, in order to normalise  $p(\theta | D)$ , we need to match  $p(\theta | D)$  w/ Beta,

$$\text{so the appropriate normalising constant} = \frac{\Gamma(|T|+a+|H|+b)}{\Gamma(|T|+a)\Gamma(|H|+b)}$$

$$\therefore p(\theta | D) = \frac{\Gamma(|T|+a+|H|+b)}{\Gamma(|T|+a)\Gamma(|H|+b)} \theta^{|T|+a-1} (1-\theta)^{|H|+b-1} = \text{Beta}(|T|+a, |H|+b)$$

With more data, the posterior dist. became more penky  $\rightarrow$  more certain about our estimate  $\theta$

### Posterior Predictive Distribution (predicting unknown data)

$$p(f | D, a, b) = \text{posterior predictive dist. } p(\theta | D, a, b) = \text{posterior dist.}$$

Using sum rule of probability

$$\begin{aligned} p(f | D, a, b) &= \int_0^1 p(f, \theta | D, a, b) d\theta \\ &= \int_0^1 p(f | \theta, D, a, b) p(\theta | D, a, b) d\theta \\ &= \int_0^1 p(f | \theta) p(\theta | D, a, b) d\theta \end{aligned}$$

if we know  $\theta$ , the next prediction  $f$  is independent of previous data  $D, a, b$

basically a weighted average for each prediction

$p(\theta | D, a, b)$  act as a weight of a prediction, how much do you believe in the  $\theta$

The prediction prob  $p(f | \theta)$  is depend on  $\theta$  and if the prob of  $\theta$  showing up is "small", then the overall predictive prob will be even lower.

MLE & MAP return  $\theta_{\text{MLE}}$  &  $\theta_{\text{MAP}}$  which can be use in the probability to get the highest probability

Fully Bayesian return full distribution of  $p(\theta | D)$

As we observe lots of data, the diff in prediction btw this 3 method become less noticeable

## Lecture 4: Linear Regression

Dataset  $D = \{(x_i, y_i)\}_{i=1}^N$        $X = \{x_1, \dots, x_N\}$  = observation     $x_i \in \mathbb{R}^D$      $X \in \mathbb{R}^{N \times D}$   
 $y = \{y_1, \dots, y_N\}$  = target     $y_i \in \mathbb{R}$   
 $y_i = f(x_i) + \varepsilon_i$      $\varepsilon_i \sim N(0, \beta^{-1})$     noise

there's always a bias term

feature vector  $\vec{x} = (1, x_1, \dots, x_D)^T$  weight vector  $\vec{w} = (w_0, w_1, \dots, w_D)^T$

Choose best  $w$  to fit the data

Loss function - measure error b/w our model (parameterized by  $w$ ) and observed data  $D = \{(x_i, y_i)\}_{i=1}^N$

↪ 1) Least Square  $\Rightarrow E_{LS}(w) = \frac{1}{2} \sum_{i=1}^N (f_w(x_i) - y_i)^2 = \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2$

Find optimal  $w^*$  to min error

$$w^* = \arg \min_w \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2$$

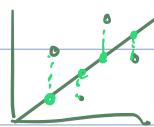
$$\text{As matrix form} \rightarrow w^* = \min_w \frac{1}{2} (Xw - y)^T (Xw - y)$$

$$\nabla_w E_{LS}(w) = \frac{1}{2} \frac{\partial}{\partial w} (w^T X^T X w - 2 w^T X^T y + y^T y)$$

$$= \frac{1}{2} (2 X^T X w - 2 X^T y) = X^T X w - X^T y = 0$$

$$w^* = \underbrace{(X^T X)^{-1}}_{X^+} X^T y = \text{Normal equation of LS problem}$$

Note  $\frac{\partial}{\partial x} x^T A x = 2Ax$   
 $\text{or } (A+A^T)x \text{ if } A \text{ is not symmetric}$



## Non-linear Dependency

Polynomial - Universal function approx.

$$\text{for 1D} \rightarrow f_w(x) = w_0 + \sum_{j=1}^M w_j x^j$$

$$\text{or more generally } f_w(x) = w_0 + \sum_{j=1}^M w_j \phi_j(x) = w^T \phi(x) \quad \text{where } \phi(x) \text{ is non linear, but the whole function } f_w(x) \text{ is still in linear form}$$

$\phi(x)$  can be e.g. 1. Polynomial -  $\phi_j(x) = x^j$  ( $x, x^2, x^3, \dots$ )

2. Gaussian -  $\phi_j(x) = \exp(-\frac{(x - \mu_j)^2}{2s^2})$

3. Sigmoid -  $\phi_j(x) = \sigma(a) = \frac{1}{1 + e^{-a}}$

use along with  $E_{LS}(w) = \frac{1}{2} \sum_{i=1}^N (w^T \phi(x_i) - y_i)^2 = \frac{1}{2} (\phi(x) w - y)^T (\phi(x) w - y)$

and  $w^* = \min_w E_{LS}(w) = \underline{(\phi(x)^T \phi(x))^{-1} \phi(x)^T y}$  (It was  $(X^T X)^{-1} X^T y$  before but since  $X$  got transform, it effect here as well)

## Choose Degree Polynomial M

- The higher the degree, the higher the variance  $\rightarrow$  overfitting (try to fit perfectly to every point)
- The lower, the higher the bias  $\rightarrow$  underfit
- Can be max  $N-1$   $N = \#$  of data
- Add Regularization like Ridge Regression (L2) to reduce  $w$  to reduce overall Loss

$$E_{\text{Ridge}}(w) = \frac{1}{2} \sum_{i=1}^N [w^T \phi(x_i) - y_i]^2 + \frac{\gamma}{2} \|w\|_2^2$$

$\gamma$  = regularization strength  
larger  $\gamma$  lead to smaller weight

## Bias / Variance Tradeoff

High bias - model too simplistic to capture the underlying pattern of data

High variance - model try to capture everything including noise which lead to high fluctuation  
- lead to overfitting

too high  $\gamma \rightarrow$  high bias

model too high capacity /  $\gamma$  too low  $\rightarrow$  high variance

Popular technique to balance b/w the 2  
select large capacity model and keep  
variance in check by choosing appropriate  $\gamma$

## Probabilistic Linear Regression

$$y_i = f_w(x_i) + \xi_i \rightarrow \xi_i \sim N(0, \beta^{-1}) \quad \text{where } \beta = \frac{1}{\sigma^2}$$

if noise ( $\xi$ ) is gaussian, then  $y_i$  follow gaussian as well

$$y_i \sim N(f_w(x_i), \beta^{-1})$$

Likelihood of a single sample  $p(y_i | f_w(x_i), \beta) = N(y_i | f_w(x_i), \beta^{-1})$

$$\text{For dataset D} \quad p(y | X, w, \beta) = \prod_{i=1}^N p(y_i | f_w(x_i), \beta)$$

Maximum Likelihood wrt.  $w, \beta$  return parameters  $w, \beta$  that will give highest probability to see  $y$

$$w_{ML}, \beta_{ML} = \arg \max_{w, \beta} p(y | X, w, \beta)$$

Use Negative log likelihood (NLL)  $-w_{ML}, \beta_{ML} = \arg \min_{w, \beta} -\ln p(y | X, w, \beta)$

let  $E(w, \beta) = \text{Maximum likelihood Error function}$

$$E_{ML}(w, \beta) = -\ln p(y | X, w, \beta) = -\ln \left[ \prod_{i=1}^N N(y_i | f_w(x_i), \beta^{-1}) \right]$$

$$\text{In the end} \Rightarrow \mathbf{w}_{ML} = \arg \min_{\mathbf{w}} E_{LS}(\mathbf{w}) = (\phi(\mathbf{x})^T \phi(\mathbf{x}))^{-1} \phi(\mathbf{x})^T \mathbf{y}$$

Which mean Maximizing the likelihood estimate is equivalent to minimizing the Least Square Error Func.

now a constant

$$\beta_{ML} = \arg \min_{\beta} E_{LS}(\mathbf{w}_{ML}^\top / \beta)$$

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}_{ML}^\top \phi(\mathbf{x}_i) - y_i)^2 = \text{Mean Square Error}$$

→ Use MLE can easily lead to overfit → use MAP & Posterior Distribution instead  
find  $\mathbf{w}, \beta$  on Posterior Distribution, MAP, Fully Bayesian

$$p(\mathbf{w} | \mathbf{x}, \mathbf{y}, \beta, \alpha) = \frac{p(\mathbf{y} | \mathbf{x}, \mathbf{w}, \beta) \cdot p(\mathbf{w} | \alpha)}{p(\mathbf{y} | \mathbf{x}, \beta, \alpha)} \quad \alpha = \text{precision of a distribution}$$

if  $\alpha \rightarrow 0 \quad \mathbf{w}_{MAP} \rightarrow \mathbf{w}_{ML}$

choose prior  $p(\mathbf{w} | \alpha)$ ?  $\Rightarrow$  Gaussian b/c likelihood is also gaussian

$$w_{MAP} = \max_{\mathbf{w}} \underbrace{\frac{p(\mathbf{y} | \mathbf{x}, \mathbf{w}, \beta)}{p(\mathbf{y} | \mathbf{x}, \beta, \alpha)} \cdot p(\mathbf{w} | \alpha)}_{\text{constant} = \text{ignore}} = \max_{\mathbf{w}} \ln p(\mathbf{y} | \mathbf{x}, \mathbf{w}, \beta) + \ln p(\mathbf{w} | \alpha) - \ln p(\mathbf{y} | \mathbf{x}, \beta, \alpha) = \arg \min_{\mathbf{w}} -\ln p(\mathbf{y} | \mathbf{x}, \mathbf{w}, \beta) - \ln p(\mathbf{w} | \alpha)$$

$$\text{let } E_{MAP} = -\ln p(\mathbf{y} | \mathbf{x}, \mathbf{w}, \beta) - \ln p(\mathbf{w} | \alpha) + \text{const}$$

$$\text{In the end } E_{MAP} \propto E_{Ridge}(\mathbf{w}) + \text{const} = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^\top \phi(\mathbf{x}_i) - y_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 + \text{const} \quad \text{where } \lambda = \frac{\alpha}{\beta}$$

$$\beta = \frac{\lambda}{\lambda}$$

Compute full posterior distribution

$$p(\mathbf{w} | \mathbf{D}) \propto p(\mathbf{y} | \mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha)$$

Since both likelihood & prior is Gaussian, the poster is gaussian as well

$$p(\mathbf{w} | \mathbf{D}) \sim N(\mathbf{w} | \mu, \Sigma) \quad \text{w/ } \mathbf{w}_{MAP} = \mu \quad (\text{mode} = \text{mean})$$

$$\text{if } \alpha \rightarrow 0 \quad \mathbf{w}_{MAP} \rightarrow \mathbf{w}_{ML}$$

Posterior Predictive Distribution - predict new data ( $\mathbf{x}_{\text{new}}, \mathbf{y}_{\text{new}}$ )

Plugging in estimated parameter  $\mathbf{w}_{ML}, \beta_{ML} / \mathbf{w}_{MAP}, \beta_{MAP}$ , we get Predictive Distribution

Recall that to make new prediction  $\hat{y}_{\text{new}}$  for  $\mathbf{x}_{\text{new}}$

$$y_i \sim N(f_w(\mathbf{x}_i), \beta^2) \text{ so } p(y_i | \mathbf{x}_i, \mathbf{w}, \beta) = N(y_i | \mathbf{w}^\top \phi(\mathbf{x}_i), \beta^2)$$

For  $y_{\text{new}}, x_{\text{new}}$

w is constant variance is fixed

$$\text{For ML : } p(y_{\text{new}} | x_{\text{new}}, w_{\text{ML}}, \beta_{\text{ML}}) = N(\hat{y}_{\text{new}} | w_{\text{ML}}^T \phi(x_{\text{new}}), \sigma_{\text{ML}}^2)$$

$$\text{For MAP : } p(\hat{y}_{\text{new}} | x_{\text{new}}, w_{\text{MAP}}, \beta) = N(\hat{y}_{\text{new}} | w_{\text{MAP}}^T \phi(x_{\text{new}}), \sigma_{\text{MAP}}^2) \quad \sigma_{\text{MAP}}^2 = \frac{\alpha}{n}$$

Alternatively, we can use the full posterior dist.  $p(w|D)$

This allow us to compute Posterior Predictive Distribution

w is now random variable and not fixed like in Predictive Distribution  $(w_{\text{ML}}, w_{\text{MAP}})$

$$p(\hat{y}_{\text{new}} | x_{\text{new}}, D) = \int p(\hat{y}_{\text{new}}, w | x_{\text{new}}, D) dw \quad \text{it has its own distribution}$$

$$= \int \underbrace{p(\hat{y}_{\text{new}} | x_{\text{new}}, w)}_{\text{Gaussian}} \cdot \underbrace{p(w | D)}_{\substack{\text{Gaussian} \\ D \text{ is gone b/c already have } w}} dw \quad x_{\text{new}} \text{ is gone b/c it's independent of } w$$

$$= N(\hat{y}_{\text{new}} | \mu^* \phi(x_{\text{new}}), \sigma^2 + \phi(x_{\text{new}})^T \Sigma \phi(x_{\text{new}}))$$

" $w_{\text{MAP}}$ "

Variance change with new data  
(uncertainty)

Advantage - More accurate estimate about uncertainty in the prediction e.g. data-dependent uncertainty estimate

## Lecture 5 : Linear Classification

Classification - Output  $y$  belong to one of  $C$  classes

Regression - Output  $y$  is continuous

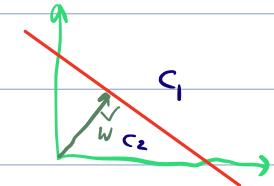
### Zero-one loss (loss function)

- denote number of misclassified samples

$$l_0(y, \hat{y}) = \sum_{i=1}^n I(\hat{y}_i \neq y_i) \rightarrow \text{sum up all the time that predicted wrong}$$

### Hyperplane as a decision boundary

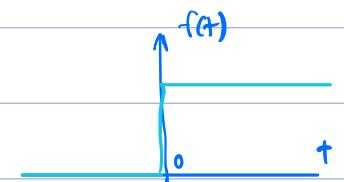
For 2 Class - Separate points from 2 classes by a hyperplane



$$\text{hyperplane} = w^T x + w_0 = \begin{cases} 0 & \text{if } x \text{ is on the plane} \\ >0 & \text{if } x \text{ is on normal side of plane} \\ <0 & \text{else} \end{cases} \quad \begin{matrix} w = \text{normal vector} \\ w_0 = \text{offset} \end{matrix}$$

### Perceptron

Decision Rule :  $y = f(w^T x + w_0)$  where  $f(t) = \begin{cases} 1 & \text{if } t > 0 \\ 0 & \text{else} \end{cases}$

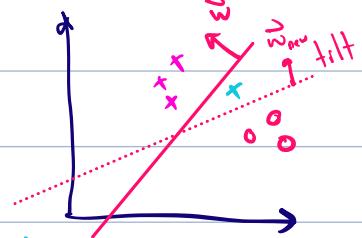


### Learning Rule for Perceptron

- tilt the hyperplane to include the  $x$  that was misclassified

$$w \leftarrow \begin{cases} w + x_i & \text{if } y_i = 1 \\ w - x_i & \text{if } y_i = 0 \end{cases} \quad w_0 = \begin{cases} w_0 + 1 & \text{if } y_i = 1 \\ w_0 - 1 & \text{if } y_i = 0 \end{cases}$$

(I think  $y_i = 0$  is for class above hyperplane)



- only work if data is linearly separable

so if it misclassified (that means  $\hat{y}=1$ )  
then the hyperplane has to be tilted down to include the point

- can scale up to multiple class by learning many hyperplane

### One-versus-rest classifier

- for  $K$  class, use  $K-1$  classifier. Each classifier solves a 2 class problem ( $C_k$  or not  $C_k$ )

- Each point is classified to class  $C$  if the distance from hyperplane is maximal

### One-versus-one classifier

$\binom{K}{2}$

(deep inside class  $C$ ) (we don't know if it's interested)

- for  $K$  class, introduce  $K(K-1)/2$  binary classifier (one for every pair of classes)

- Each classifier classifies point to either  $C_i$  or  $C_j$

- Each point is classified to class  $C$  according to majority vote amongst all classifier (we don't know if the point get some amount of vote)

- Multiclass Discriminant - train each classifier ( $C$  linear functions)

$$f_c(x) = w_c^T x + w_{0c} \quad (\text{hyperplane}) \quad \text{w/ decision rule } \hat{y} = \arg \max_{c \in C} f_c(x)$$

Divide the region, not based on side of hyperplane, but based on the distance predict class for point  $x$

from the boundary. (The decision boundary is where  $y_k(x) = y_j(x) \forall j \neq k$  w/ the class that has the

If point lies on the decision boundary, it can either be  $k$  or  $j$ . It's undecided) biggest  $f_c(x)$

## Basis functions

Apply a non-linear transformation  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$  (can learn the transformation w/ NN)

linear transformation can stretch, rotate or shift but cannot separate data and make it linearly separable.

## Limitation of Hard-decision based classifier

- 1) No measure of uncertainty (probability for each class)
- 2) Can't handle noisy data
- 3) Poor generalisation (try too hard to classify all points correctly)  $\rightarrow$  difficult to optimise

## Probabilistic models for Classification

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

- 2 Type of model
  - Generative  $\rightarrow$  Model joint distribution (get  $p(y=c|x)$  by finding  $p(y=c) \& p(x|y=c)$  first)
  - Discriminative  $\rightarrow$  Directly model the distribution  $p(y=c|x)$

Given  $p(y=c|x)$  we can make prediction  $\hat{y}$   $\hat{y} = \max_{c \in C} p(y=c|x)$

## Generative models for linear classification

$$p(y=c|x) \propto p(x|y=c)p(y=c)$$

$p(x|y=c)$  - class conditions - probability of generating point  $x$  given that it belongs to class  $c$

$p(y=c)$  - class prior - prior prob. of a point belonging to class  $c$

## Applying generative model

- Choose parametric model for class conditional and class prior (Bernoulli if  $\text{class}=2$  or  $\text{+} \sim \text{Normal}$  if  $\text{class}>2$ )

$$p(x|y=c, \psi) \text{ and } p(y=c|\theta)$$

estimate model's parameter  $\{\psi, \theta\}$  from data  $D$  (using e.g. ML) - obtain estimate  $\{\hat{\psi}, \hat{\theta}\} \rightarrow$  this step is called learning

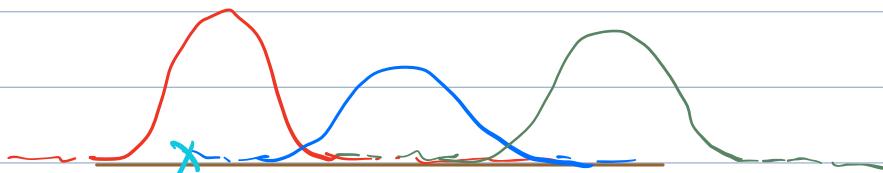
Once fitted, we can perform inference - classify label for  $x$  using Bayes rule

$$p(y=c|x, \hat{\psi}, \hat{\theta}) \propto p(x|y=c, \hat{\psi}) p(y=c|\hat{\theta})$$

Can also generate new data

→ sample a class label  $y_{new} \sim p(y|\hat{\theta})$  (predicting it's a cat)

→ sample new data  $x_{new} \sim p(x|y=y_{new}, \hat{\psi})$  (predicting value of pixel  $x$  knowing it's a cat)



let  $\begin{bmatrix} r \\ g \\ b \end{bmatrix} = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.4 \end{bmatrix} = p(y|\hat{\theta})$  prob of each class

$x_{new} \rightarrow p(y=\text{red}|x, \hat{\psi}, \hat{\theta}) = p(x|y=\text{red}, \hat{\psi}, \hat{\theta}) p(y=\text{red}) = \underline{\text{big prob}} \times 0.3 = \underline{\text{big prob}}$

$$\rightarrow p(y=\text{blue}|x, \hat{\psi}, \hat{\theta}) = \text{small} \times 0.3 = \text{small prob}$$

$$\rightarrow p(y=\text{green}|x, \hat{\psi}, \hat{\theta}) = \text{very small} \times 0.4 = \text{small prob}$$

→ after finding prob of all class( $y$ ) for point  $x$ , choose the class w/ highest prob

How to choose class prior

$p(y=c)$   $y \sim \text{Categorical}(\theta)$   $\theta \in \mathbb{R}^c$  is prob of each class

$$\text{so } p(y=c) = \theta_c \text{ and } 0 \leq \theta_c \leq 1 \text{ and } \sum_{c=1}^C \theta_c = 1$$

the Maximum Likelihood estimate for  $\theta$

$$\hat{\theta}_c^{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N I(y_i=c) \quad \text{e.g. if } b=3 \quad r=5 \quad g=2 \\ \hat{\theta}_b = \frac{3}{10} \quad \hat{\theta}_r = \frac{5}{10} \quad \hat{\theta}_g = \frac{2}{10}$$

Choose class conditional dist.

→ Use multivariate normal for each class  $p_{\text{class}}(x|y=c) = N(x|\mu_c, \Sigma)$  has  $D \times D$  param

Linear Discriminant Analysis (assume shared covariance  $\Sigma$ )

let class conditional =  $p(x|y=c) = N(x|\mu_c, \Sigma)$   $\mu_c$  = mean for each class  $\Sigma$  = same covariance matrix for all class

For  $C=2$  - Hyperplane = boundary (normally not the case)

$$p(y=1|x) = \sigma(w^T x + w_0) \xrightarrow{\text{so}} p(y=0|x) = 1 - \sigma(w^T x + w_0) \quad \text{or} \quad y \sim \text{Bernoulli}(\sigma(w^T x + w_0))$$

For  $C > 2$  - use softmax function

$$p(y=c|x) = \frac{p(x|y=c)p(y=c)}{\sum_{i=1}^C p(x|y=i)p(y=i)} = \frac{\exp(w_c^T x + w_{c0})}{\sum_{i=1}^C \exp(w_i^T x + w_{i0})}$$

LDA will learn linear decision boundary (it use hyperplane which create decision boundary)

Naive Bayes (Assume d features of sample are independent given the class)

let class conditional  $= p(x|y=c) = N(x|\mu_c, \Sigma_c) \rightarrow \mu$  and  $\Sigma$  for each class c

$$p(x_1, x_2, \dots, x_d | y=c) = \prod_{i=1}^d p(x_i | y=c)$$

Since  $x_1, \dots, x_d$  is indep. given class c, the covariance  $\Sigma_c$  is diagonal matrix  $\rightarrow$  only have D param instead of  $D \times D$

$$\text{LDA param} = C \cdot D + D \cdot P \quad \text{NB param} = C(D+D)$$

$$p(y=c|x) = \sigma(x^T w_c + w_{c0}) \rightarrow \text{quadratic func!} \rightarrow \text{quadratic decision boundary}$$

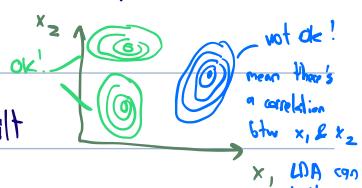
\* If  $\Sigma_c$  is same for all class  $\rightarrow$  same solution as LDA

## Pros of NB

- Allow diff covariance structure

use Gaussian distribution

- Independent feature allows to easily handle diff data types/mix data types ( $x_1$ -numerical,  $x_2$ -categorical...)



## Cons of NB

- b/c features are indep. It ignore the correlation btw the feature  $\rightarrow$  so boundary cannot be tilt

## Discriminative Model for Linear Classification (model $p(y|x)$ )

- In Generative model  $w, w_0$  are precompute before plugging into  $\sigma$  to find  $p(y=c|x) \propto p(x|y=c) p(y=c)$

$w, w_0$  also depend on parameters of class conditionals  $\mu_0, \mu_1, \Sigma$

In Discriminative model  $w$  and  $w_0$  is random variable that can be optimize

## Logistic Regression

$$\text{let } w^T x = w^T x + w_0$$

$$y|x \sim \text{Bernoulli}(\sigma(w^T x + w_0))$$

Learning logistic regression  $\rightarrow$  find good parameter  $w$  that "explain" the training set D

$$\begin{aligned} \text{(Assume sample are i.i.d.) } p(y|w, X) &= \prod_{i=1}^N p(y_i|x_i, w) \rightarrow p(y_i|x_i, w) = \text{Bernoulli}(\dots) \\ &= \prod_{i=1}^N \sigma(w^T x_i)^{y_i} (1 - \sigma(w^T x_i))^{1-y_i} \end{aligned}$$

Optimize  $w$   $\xrightarrow{\text{learning}}$  Find MLE for  $w \rightarrow w = \arg \max_w p(y|w|x)$

## Negative Log Likelihood (NLL) - Error Function

$$E(w) = -\log p(y|w, X) = -\sum_{i=1}^n (y_i \log \sigma(w^T x_i) + (1-y_i) \log (1-\sigma(w^T x_i))) \rightarrow \text{this loss function is Binary Cross Entropy}$$

↳ no  $\frac{1}{n}$  b/c use total loss in ML lecture

and Finding MLE for  $w$  is equivalent to minimizing  $E(w)$

$$w^* = \arg \min_w E(w)$$

But b/c MLE can lead to overfitting, we can use MAP estimation which will include regularization

$$E(w) = -\log p(y|w, X) + \gamma \|w\|_2^2 \text{ (for L-2 norm / ridge)}$$

## Multiclass Logistic Regression

Use Softmax  $p(y=c|x) = \frac{\exp(w_c^T x)}{\sum_{i=1}^c \exp(w_i^T x)}$

And optimal  $w \rightarrow w^* = \arg \max_w p(Y|w, X) \rightarrow$  use NLL  $\rightarrow -\log p(Y|w, X) = E(w)$

$$E(w) = -\sum_{i=1}^n \sum_{c=1}^C y_{ic} \log \frac{\exp(w_c^T x_i)}{\sum_{j=1}^C \exp(w_j^T x_i)} \text{ where } y_{ic} = \begin{cases} 1 & \text{if sample } i \text{ belongs to class } c (i=c) \\ 0 & \text{ow.} \end{cases}$$

This is called Cross Entropy ↳ no  $\frac{1}{n}$  b/c use total loss in ML lecture

## Generative VS Discriminative Model

- Discriminative model achieve better perf in general for classification task (less assumption required)
- Generative work well but fragile when assumption are violate (same covariance / independent feature)
- Generative models provide benefits of better handling missing data, detecting outliers, generating new data  
↳ create new data to fill in the blanks ↳ if data has low prob in all class, it's most likely an outliers

## Note to Recap

Linear Classification for C-2 → use Bernoulli to predict output (label)

→ Find optimal  $w^* \rightarrow \arg \max_w p(y|w, x) \quad y | x \sim \text{Bernoulli}$

→ Use NLL → get loss function  $E(w) = \text{BCE} = -\sum y \log y + (1-y) \log (1-y)$

→  $w^* = \arg \max_w p(y|w, x) = \min_w E(w)$  (no closed form)

Linear Regression → use Gaussian to predict output (continuous)

→ Find optimal  $w^* \rightarrow \arg \max_{w, \beta} p(y|w, x, \beta) = \arg \max_{w, \beta} N(f_w(x), \beta^{-1})$

→ Use NLL → get loss function  $E_{LS}(w) = \text{Least square} = \frac{1}{2} (f_w(x) - y)^2$

→  $w^* = \arg \max_{w, \beta} p(y|w, x, \beta) = \min_w E_{LS}(w) = (X^T X)^{-1} X^T y$

## Lecture 6: Optimization (Convex!)

Denote  $f(\theta) = \text{objective function}$

Find minimum of objective function  $\rightarrow$  convex

### Convexity

$X$  is convex set iff  $\forall x, y \in X \quad \exists_{\lambda} \in (0, 1) \quad \lambda x + (1-\lambda)y \in X$

$f(x)$  is convex function iff  $\forall x, y \in X \quad f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$

For convex func - each local minimum is a global minimum (could be more than 1 global min)

- if  $f_\theta$  is convex,  $\nabla f_\theta(\theta^*) = 0$  then  $\theta^*$  is a global minimum

- Rate of rise on function  $= f(y) - f(x) \geq \frac{f((1-\lambda)x + \lambda y) - f(x)}{\lambda}$  or  $f(y) - f(x) \geq (y-x)^T \nabla f(x)$  when  $\lambda \rightarrow 0$

$\therefore$  if  $f: X \rightarrow \mathbb{R}$  is differentiable and  $X$  is convex, then  $f$  is convex iff  $\forall x, y \in X \quad f(y) \geq f(x) + (y-x)^T \nabla f(x)$

### Verify Convexity

1)  $\forall f(x) + (1-\lambda)f(y) \geq \underset{\text{convex}}{f(\lambda x + (1-\lambda)y)}$  for all  $x, y \in X$

2) Hessian matrix is PSD on the set  $(x^T \nabla^2 f(x, \dots), x)$  then func is convex on the convex set

(for 1 variable it's convex on a domain, if  $\nabla^2 f(x) \geq 0$  for all  $x$  in the domain)

if  $\nabla^2 f(x)$  or  $x^T H(x, \dots) x \leq 0$ , then func is concave

Linear func is both convex and concave

3) Obtain by operation that preserve convexity

Let  $f_1$  and  $f_2$  be convex func and let  $g$  be concave

-  $h(x) = f_1(x) + f_2(x)$  is convex

-  $h(x) = \max \{f_1(x), f_2(x)\}$  is convex

-  $h(x) = c \cdot f_1(x)$  is convex if  $c \geq 0$

-  $h(x) = c \cdot g(x)$  is convex if  $c \leq 0$

-  $h(x) = f_1(Ax+b)$  is convex ( $A$  matrix,  $b$  vector)

-  $h(x) = m(f_1(x))$  is convex if  $m: \mathbb{R} \rightarrow \mathbb{R}$  is convex and non-decreasing

let  $z(x) = h(f(x)) \rightarrow z'(x) = h'(f(x)) \cdot f'(x) \rightarrow z''(x) = h''(f(x)) \cdot f'(x) \cdot f'(x) + f''(x) \cdot h'(f(x))$

Proof

and  $\gamma$  is convex

$$= h'(f(x)) + f'(x)h(f(x)) \geq 0$$

must be  $\geq 0$        $\geq 0$        $\geq 0$        $\downarrow$   
 so  $h$  must be convex      b/c  $f(x)$  is convex (non-decreasing)

## Convex Set

1. Prove definition  $\forall x + (1-\gamma)y$

2. if  $A$  and  $B$  are convex set, then  $A \cap B$  is convex set

$$y = p_0 + p_1 x + \varepsilon \quad (\text{line})$$

Unconstrained Minimization problem  $\rightarrow$  e.g. Ordinary Least Square Regression  $\rightarrow \theta^* = f'(\theta) = 0$

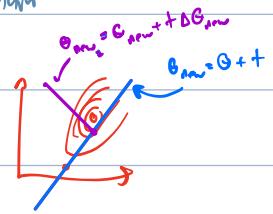
## First Order Optimization (Exploit gradient)

### - Gradient Descent

Gradient - point at steepest ascent direction

GD w/ Line Search - get descent direction, then  $\theta_{\text{new}} = \theta + t \Delta \theta$

- Turn multidimension prob into one-dimension prob (do one dim at a time)



Step 1)  $\Delta \theta = -\nabla f(\theta)$     2) Line Search:  $t^* = \min_{t>0} f(\theta + t \Delta \theta)$     3) Update:  $\theta_{\text{new}} = \theta + t^* \Delta \theta$

$\nabla f(\theta)$  doesn't point to global min, it point to local min.

You can run GD Line Search in parallel (partial gradient) and aggregate them.

But line search is expensive, for each tested step size  $\rightarrow$  scan through all datapoints

### - Learning Rate

Instead of using line search to find optimal  $t \rightarrow$  simply pick a learning rate  $\gamma$

$$\theta_{+1} = \theta_+ - \gamma \nabla f(\theta_+) \quad \gamma \rightarrow \text{too small} \rightarrow \text{slow convergence, end up in local min or saddle points} \quad \gamma \rightarrow \text{too large} \rightarrow \text{oscillate, never converge}$$

Solution - Learning rate schedule  $\rightarrow$  Let  $\gamma$  be decreasing function of iteration  $t$

Momentum -  $\gamma$  increase (accelerate) as long as gradients keep pointing at the same direction

$$\theta_{+1} \leftarrow \theta_+ - m_+ \quad m_+ \leftarrow \gamma \nabla f(\theta_+) + \beta m_{+1} \quad \beta = \text{momentum factor}$$

Adagrad - Diff learning rate per parameter

$$\begin{aligned} \text{Adam} - m_+ &= \beta_1 m_{+1} + (1-\beta_1) \nabla f(\theta_+) \quad m_+ = \frac{m_+}{1-\beta_1} \quad (\text{similar to momentum}) \quad \text{gradient fluctuate a lot} \\ - v_+ &= \beta_2 v_{+1} + (1-\beta_2) \nabla f(\theta_+)^2 \quad v_+ = \frac{v_+}{1-\beta_2} \quad (\text{if gradient is large, it will decrease the step}) \\ - \theta_{+1} &= \theta_+ - \frac{\gamma}{\sqrt{v_+ + \varepsilon}} \hat{m}_+ \quad \text{if } v_+ \text{ small, } \rightarrow \text{increase } \gamma \end{aligned}$$

## Second Order - Newton

$\nabla^2 f(\theta) \succ 0$  = Hessian is PSD

↓  
step

Taylor Expansion of  $f$  at point  $\theta_+$   $f(\theta_+ + \delta) = f(\theta_+) + \delta^\top \nabla f(\theta_+) + \frac{1}{2} \delta^\top \nabla^2 f(\theta_+) \delta + O(\delta^3)$  :=  $g(\delta)$

Minimize Approx. leads to  $\nabla g(\delta) = \nabla f(\theta) + \nabla^2 f(\theta) \delta = 0 \rightarrow \delta = (\nabla^2 f(\theta))^{-1} \nabla f(\theta)$

And the update step  $\theta_{t+1} = \theta_t - [\nabla^2 f(\theta_t)]^{-1} \nabla f(\theta_t)$  Repeat until convergence

+ Fast convergence (require less iteration) - Require compute inverse Hessian ( $O(d^3)$ )  $\Rightarrow$  Only use for low dim problem

## SGD

(Exact) expectation can be approx. by smaller sample ( $S \subseteq \{1, \dots, n\}$ )

Instead of computing the loss function for all instances, only compute loss func to the subset of data.

$$\mathbb{E}_{i \in \{1, \dots, n\}} [L_i(\theta)] \approx \frac{1}{|S|} \sum_{j \in S} L_j(\theta) \rightarrow \frac{1}{n} \sum_{i=1}^n L_i(\theta) \approx \frac{1}{|S|} \sum_{j \in S} L_j(\theta) \rightarrow \sum_{i=1}^n L_i(\theta) \approx \frac{n}{|S|} \sum_{j \in S} L_j(\theta)$$

original expectation loss      batch mean loss

times  $n$  to upscale  
 to original loss  
 $n = \#$  of whole data

The gradient computation can also be done using only subset of data  $\rightarrow$  noisy unbiased estimate

1. Randomly pick small subset  $S$  of the data (mini-batch)
2. Compute gradient based on mini-batch
3. Update weight  $\theta_{t+1} \leftarrow \theta_t - \gamma \left( \frac{1}{|S|} \sum_{j \in S} \nabla L_j(\theta_t) \right)$  (Update happen for every minibatch/iteration)
4. Pick new subset and repeat 2-4

Larger minibatch lead to a more stable gradient (smaller variance in the estimated gradient)

SGD converge almost surely to global min when objective func is convex (and local min for non-convex)

Converge need more iteration to converge than GD but b/c it use minibatch, each iteration is faster

## Distributed Learning

- Distribute computation across multiple machine

Data Parallelism (CNN) vs Model Parallelism (SGD)

## Lecture 7/8: Deep Learning

Basis function  $\phi(x)$  - Apply (non-linear) transformation  $\phi$  that map samples to a space where they are linearly separable

Multi-layer Perceptron (MLP) = Fully Connected NN Feed forward NN = NN w/ 1 hidden layer

Most activation func apply element-wise. Exception is softmax (use all input for normalize)

Use activation func to allow model to solve nonlinear problem.

## Universal Approximation Theorem

MLP w/ linear output layer and 1 hidden layer can approx. any continuous function defined over a closed and bounded subset of  $\mathbb{R}^D$  with constraint of activation func and given #of hidden units are large enough  
 In short: there exist MLP that can approx. any func well

## Loss function

Diff task require changing

1. Activation function in the final layer

2. Loss function

Prediction target	$p(y x)$	Final layer Acti. func.	Loss Func
Binary ( $C=2$ )	Bernoulli	Sigmoid	BCE
Discrete ( $C > 2$ )	Categorical	Softmax	CE
Continuous	Gaussian	Identity $\Leftrightarrow$ (basically identity)	MSE (Mean Square Error)
Other loss func			$E(w) = -\frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2$

- Mean Absolute Error - useful if we have outliers
- Huber Loss - MSE good at close dist, MAE good at far away point, this func combine both benefit

## Backpropagation

1. Write func as composition of modules

2. Work out local derivative of each module symbolically

3. Do forward pass for given input  $x \rightarrow$  cache the result for each func. to be use in backprop

4. Compute local derivative for  $x$

5. Obtain global derivative by multiplying local deriv.

$\textcircled{A} \rightarrow \textcircled{B} \rightarrow \textcircled{C}$  Compute deriv.

$\textcircled{A} \rightarrow \textcircled{B} \rightarrow \textcircled{C}$  along each paths

$$c(a,b) = a(x) + b(y)$$

$$c(a,b) = a(x) + b(y)$$

$$\frac{\partial c}{\partial x} = \frac{\partial c}{\partial a} \cdot \frac{\partial a}{\partial x} + \frac{\partial c}{\partial b} \cdot \frac{\partial b}{\partial x}$$

$$\frac{\partial c}{\partial y} = \frac{\partial c}{\partial a} \cdot \frac{\partial a}{\partial y} + \frac{\partial c}{\partial b} \cdot \frac{\partial b}{\partial y}$$

$$\textcircled{A} \rightarrow \textcircled{B} \rightarrow \textcircled{C} \quad \frac{\partial c}{\partial x} = \frac{\partial c}{\partial a} \cdot \frac{\partial a}{\partial x} + \frac{\partial c}{\partial b} \cdot \frac{\partial b}{\partial x}$$

Multivariate Chain rule

$$\frac{\partial c}{\partial x} = \sum_{i=1}^m \frac{\partial c}{\partial a_i} \cdot \frac{\partial a_i}{\partial x}$$

## Jacobian and gradient

Consider func  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  (take  $n$  variable, output  $m$  var)

or func.

$$f(x_1, x_2, \dots, x_n)$$

$$\text{let } a = f(x)$$

$$J = \frac{\partial \mathbf{a}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial a_1}{\partial x_1} & \dots & \frac{\partial a_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial a_m}{\partial x_1} & \dots & \frac{\partial a_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

es. loss func.  
let  $g: \mathbb{R}^m \rightarrow \mathbb{R}$  and let  $c = g(a)$   
 $\downarrow g(a_1, \dots, a_m)$

the gradient  $\nabla_a c$  is the transpose of the Jacobian  $\frac{\partial c}{\partial a} \in \mathbb{R}^{l \times m}$

$$c = g(a_1, \dots, a_m) \quad \frac{\partial c}{\partial a} = \left( \frac{\partial c}{\partial a_1}, \dots, \frac{\partial c}{\partial a_m} \right)$$

$$\nabla_a c = \left( \frac{\partial c}{\partial a} \right)^T = \begin{pmatrix} \frac{\partial c}{\partial a_1} \\ \vdots \\ \frac{\partial c}{\partial a_m} \end{pmatrix} = \left( \frac{\partial c}{\partial a_1}, \dots, \frac{\partial c}{\partial a_m} \right)^T$$

$\nabla_a c$  has same shape as  $a \in \mathbb{R}^{m \times l}$

and  $\frac{\partial c}{\partial x_j} = \sum_{i=1}^m \frac{\partial c}{\partial a_i} \cdot \frac{\partial a_i}{\partial x_j}$  can be write in matrix form  $\frac{\partial c}{\partial x} = \frac{\partial c}{\partial a} \cdot \frac{\partial a}{\partial x}$  where  $\left[ \frac{\partial a}{\partial x} \right]_{ij} = \frac{\partial a_i}{\partial x_j}$

$$\text{Also } \nabla_x c = \left( \frac{\partial c}{\partial x} \cdot \frac{\partial a}{\partial x} \right)^T = \left( \frac{\partial a}{\partial x} \right)^T \left( \frac{\partial c}{\partial a} \right)^T = \left( \frac{\partial a}{\partial x} \right)^T \nabla_a c$$

## Matrix Calculus

		Scalar	Vector	Matrix
		Scalar	Vector	Matrix
Output	Vector		Matrix	3-way tensor
	matrix			4-way tensor

$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial a} \cdot \frac{\partial a}{\partial w}$  materialise / computing  $\frac{\partial a}{\partial w}$  independently is expensive b/c it's a 3way tensor

However, looking at the  $\frac{\partial E}{\partial w} = \frac{\partial E}{\partial a} \cdot \frac{\partial a}{\partial w}$ , you can see that  $\frac{\partial E}{\partial w} = \frac{\partial E}{\partial a} \otimes \mathbf{X}$

$$\frac{\partial E}{\partial w} = \left( \frac{\partial E}{\partial a} \right)^T \mathbf{X}^T$$

Matrix operations are much more efficient thanks to vectorization (GPU)

## Convolution - Do feature extraction

Filter have  $L \times H \times C_{in} \times C_{out}$   $C_{out} = \# \text{ of filter}$   $C_{in} = \# \text{ of input's channel}$

Weight sharing - Weight of the kernel (filter) are shared for every patch

Padding - Valid - Do not use padding - reduce size in output  $P=0$

SAME - Adding padding ( $P=k/2$ ) to preserve size

$K = \text{size of filter}$

FULL - Increase output size ( $P=k-1$ )

Stride - distance btw positions the kernel is applied.  $\rightarrow S>1$  are a way to downsample the signal

The output size is  $D_{l+1} = \frac{D_l + 2P - K}{S} + 1$

Pooling - Doing aggregation & downsample at a same time  $\rightarrow$  Max pooling, mean pooling,  $L_p$  norm

## Weight Initialization

Issue of naive weight init

$$\begin{array}{ll} \text{Diagram: } & z_1 = x_1 w + x_2 w \\ & \frac{\partial z_1}{\partial w} = x_1 + x_2 \\ \text{Diagram: } & z_2 = x_1 w + x_2 w \\ & \frac{\partial z_2}{\partial w} = x_1 + x_2 \end{array}$$

1. Weight symmetry - If 2 hidden node have exactly the same bias and (initial weight w) same value incoming & outgoing weight, they will always get exactly the same gradient

- So they can never learn to extract diff' feature  
(update will lead to the same result)
- Break symmetry by initializing the weights to have small random value

## 2. Weight Scale

- If a hidden unit has a big # of input signal, small change on many of its incoming weight can cause learning to overshoot.
- If large # of output signal, can cause the same during backward pass
- Wrong scale can lead to vanishing/exploding gradient

Vanishing & Exploding gradient - Deep network can suffer from this problem

$\downarrow$  Solution  $\rightarrow$  Change architecture (batch norm, ReLU), Gradient clipping (limit gradient)

## Xavier Initialization

Perserve the signals' properties by using weight matrices w/  $\mu=0$  and

$$Var(w) = \frac{2}{\# \text{ of incoming} + \# \text{ of outgoing signal}}$$

Initial the weight for every FC layer

## Regularization

$$\min_w L + \|w\|_F$$

- Reduce Overfitting  $\hookrightarrow L_2$  norm shrink parameter equally,  $L_1$  norm to promote sparsity

Combine w/ other method like Data augmentation, injecting noise, dropout to reduce overfitting

Dropout - randomly drop some % of neuron (make input = 0) random new set for each iteration

## Batch Normalisation

- Stabilize the distribution of each layer's activation  $\rightarrow$  More stable gradient  $\rightarrow$  training is more stable
- Done after each layer  $\rightarrow$  allow deeper network

$$\hat{x} = \frac{x - \mu_B(x)}{\sigma_B(x) + \epsilon}$$

calculate  $\mu$  and  $\sigma$  for each minibatch to normalise the batch training data

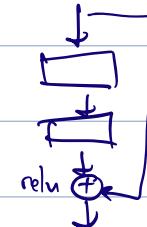
Introduce  $\gamma$  and  $\beta$  hyperparameter allow model to learn to undo the BN

$$y = \gamma \hat{x} + \beta$$

## Skip Connections

Highway Network  $\rightarrow$  gating mechanism

Residual Connections  $\rightarrow$  help w/ vanishing gradient, ignore the layer



## Lecture 9: SVM and Kernels

Linear Classifier assign class of  $x$

$$h(x) = \text{sign}(w^T x + b) \quad w \mid \text{sign}(z) = \begin{cases} -1 & \text{if } z < 0 \\ 0 & \text{if } z = 0 \\ 1 & \text{if } z > 0 \end{cases}$$

Idea: Find the hyperplane that separate data the best  $\rightarrow$  have largest margin possible  
so that new data is more likely to fall on the correct side of boundary

So we try to maximize size of margin  $m = \frac{2}{\|w\|_2 \sqrt{w^T w}}$  w/ constraint  $y_i(w^T x_i + b) \geq 1$  for all  $i$

Maximize  $\frac{2}{\|w\|_2} = \text{minimizing } \frac{1}{2} w^T w \rightarrow$  Constrained Convex Optimization Problem  $\rightarrow$  Lagrange

Lagrange dual function  $-g(\alpha) = \min_w L(w, \alpha) = \min_w (f_0(w) + \sum_i \alpha_i f_i(w))$  give lower bound on the optimal value  $f_0(\theta^*)$

Lagrange dual problem find the  $\alpha$  that gives best lower bound  $\max_{\alpha} g(\alpha)$   $\alpha \geq 0$

Apply Lagrange to SVM -  $\underbrace{y_i(w^T x_i + b) - 1}_{\geq 0} \rightarrow -f_i(y_i(w^T x_i + b) - 1) \leq 0$

$$L(w, b, \alpha) = \frac{1}{2} w^T w - \sum_i \alpha_i (y_i(w^T x_i + b) - 1)$$

(\*)  $\nabla_w L = w - \sum_i \alpha_i y_i x_i = 0 \rightarrow w = \sum_i \alpha_i y_i x_i$  weights are linear comb of training sample

$$(**) \frac{\partial L}{\partial b} = \sum_i \alpha_i y_i = 0$$

Substitute (\*) and (\*\*) back into  $g(\alpha)$  and find Lagrange Dual problem to get  $\alpha^*$

$w = \sum_i \alpha_i^* y_i x_i$  From the complementary slackness condition  $\alpha_i^* f_i(\theta^*) = 0$  (Strong Duality  $g(\alpha^*) = f_0(\theta^*)$ )

if  $\alpha_i^* \neq 0$ , then  $f_i(\theta^*) = 0 \rightarrow y_i(w^T x_i + b) - 1 = 0 \rightarrow b = y_i - w^T x_i$  (recover bias)

Furthermore,  $y_i(w^T x_i + b) - 1 = 0$  is the point at the margin

if  $\alpha_i^* = 0$ , then point is outside of margin and  $w \cdot \sum \alpha_i y_i x_i = 0$  for those point

So only the point on the margin (where  $x_i$  is support vector) are count toward min. problem (w)  $w = \sum_s \alpha_s y_s x_s$

Only need to remember few training sample  $x_i$  w/  $\alpha_i > 0$  where s = support vector point

## Soft Margin SVM

Hard Margin - doesn't work w/ data that is not linearly separable

Idea : Relax the constraint but punish the relax of constraint by introducing slack variable  $\varepsilon_i \geq 0$

for every training sample.  $\varepsilon_i$  gives the distance of how far the margin is violated in unit of  $\|w\|$

$\varepsilon_i = 0 \rightarrow$  outside of margin, correct side.  $0 < \varepsilon_i < 1$  correct side but inside margin gap  $\varepsilon_i \geq 1$  - wrong side

So the new constrain func is  $y_i(w^T x_i + b) \geq 1 - \varepsilon_i$  for all i and  $\varepsilon_i \geq 0$  (also include margin gap but wrong side)

And the new cost func is  $f_0(w, b, \varepsilon) = \frac{1}{2} w^T w + C \sum_{i=1}^N \varepsilon_i$  where  $C =$  how heavy a violation is punish.  
if  $C \rightarrow \infty$ , then we have hard margin SVM

$$L(w, b, \varepsilon, \alpha, \mu) = \frac{1}{2} w^T w + C \sum_{i=1}^N \varepsilon_i - \sum_i \alpha_i (y_i(w^T x_i + b) - 1 + \varepsilon_i) - \sum_i \mu_i \varepsilon_i$$

Lagrange variable

$$\nabla_w L = w - \sum \alpha_i y_i x_i = 0 \quad \frac{\partial L}{\partial b} = \sum_i \alpha_i y_i = 0 \quad \frac{\partial L}{\partial \varepsilon} = C - \alpha_i - \mu_i = 0 \rightarrow \alpha_i = C - \mu_i \quad \begin{matrix} \text{since } \alpha_i \geq 0 \text{ and } \mu_i \geq 0 \\ \text{then } 0 \leq \alpha_i \leq C \end{matrix}$$

Substitute into  $g(\alpha)$  and find maximum  $\alpha^*$  with Dual Problem and get same solution as Hard margin except  $0 \leq \alpha_i \leq C$  so  $\alpha_i$  cannot go to infinity.

$$\varepsilon_i \geq 1 - y_i(w^T x_i + b)$$

Look at constrain func -  $y_i(w^T x_i + b) \geq 1 - \varepsilon_i$  and  $\varepsilon_i \geq 0 \rightarrow$  we can rewrite into unconstrain opt. problem.

if point is on the wrong side / on the margin  $y_i(w^T x_i + b) = 1 - \varepsilon_i$  then  $\varepsilon_i = 1 - y_i(w^T x_i + b)$

but b/c for wrong side  $y_i(w^T x_i + b) < 1 \rightarrow \varepsilon_i > 0 = 1 - y_i(w^T x_i + b)$

if point is on the correct side and outside margin  $y_i(w^T x_i + b) = 1 - \varepsilon_i$  then  $\varepsilon_i = 1 - y_i(w^T x_i + b)$

b/c for correct side  $y_i(w^T x_i + b) \geq 1 \rightarrow \varepsilon_i \leq 0$  but b/c  $\varepsilon_i \geq 0 \rightarrow \varepsilon_i = 0$

$$\varepsilon_i = \begin{cases} 1 - y_i(w^T x_i + b) & \text{if } y_i(w^T x_i + b) < 1 \\ 0 & \text{else} \end{cases} \quad \text{Hinge loss} = \max(0, 1 - y_i(w^T x_i + b)) \rightarrow 0 = \text{predict correct} \\ 1 = \text{misclassified}$$

And we can rewrite objective func as an unconstrained optimization problem

$$\min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^N \max\{0, 1 - y_i(w^T x_i + b)\}$$

We can optimize hinge loss directly using gradient based method

## Kernels

Kernel function  $k: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$   $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$

Kernel represent inner product in the feature space spanned by  $\phi$ , this make our model nonlinear

Why use kernel instead of basis function  $\phi$

- Computing basis function can be too expensive, but directly evaluating kernel is often easier than computing  $\phi(x)$
- Kernel allow us to compare similarity btw arbitrary data (graph, non-numerical data)

## Kernel preserving operation

let  $k_1: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$   $k_2: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  be kernel

$$K(x_1, x_2) = k_1(x_1, x_2) + k_2(x_1, x_2)$$

$$K(x_1, x_2) = C \cdot k(x_1, x_2) \quad \text{if } C > 0$$

$$K(x_1, x_2) = K_1(x_1, x_2) \cdot K_2(x_1, x_2)$$

$$K(x_1, x_2) = K_3(\phi(x_1), \phi(x_2)) \quad \text{where } K_3: \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R} \text{ and } \phi: X \rightarrow \mathbb{R}^M$$

$$= \phi(x_1)^T \phi(x_2) \rightarrow \text{eg. Prove that 1 is kernel} \rightarrow K(x_1, x_2) = k_3(\phi(x_1), \phi(x_2)) = \phi(x_1)^T \phi(x_2)$$

$$\text{let } \phi(z) = 1 \text{ so } K(x_1, x_2) = (1)^T (1) = 1$$

$$\text{eg. } \exp\left(-\frac{x_1^T x_2}{2\sigma^2}\right) \exp\left(-\frac{x_2^T x_2}{2\sigma^2}\right) \rightarrow \text{let } \phi(z) = \exp\left(-\frac{z^T z}{2\sigma^2}\right)$$

$$K(x_1, x_2) = x_1^T A x_2 \quad \text{w/ } A \in \mathbb{R}^{N \times N} \text{ and } A \text{ is symmetric and PSD}$$

$$K_3(\phi(x_1), \phi(x_2)) = \underbrace{\exp\left(-\frac{x_1^T x_2}{2\sigma^2}\right)^T}_{\text{scale so}} \exp\left(-\frac{x_2^T x_2}{2\sigma^2}\right) \quad \text{Transpose doesn't matter}$$

Kernel - Polynomial =  $K(a, b) = (a^T b)^p$  or  $(a^T b + 1)$

$$\text{- Gaussian Kernel} = K(a, b) = \exp\left(-\frac{\|a-b\|^2}{2\sigma^2}\right)$$

if 2 point are similar to each other, then  $K$  is large. If  $\sigma$  is smaller, the decision boundary is more fine grain  $\rightarrow$  overfit  
look at smaller neighborhood area

## Classify new point w/ Kernelised SVM

Due to complementary slackness ( $g(x^*) = f_\theta(x^*)$  - strong duality)

This mean that both constrain of Soft margin SVM must be 0

$$(*) \alpha_i^* f_i(\theta^*) = 0 \quad \text{and} \quad (*) \mu \varepsilon_i = 0 \quad (\alpha = C - \mu)$$

Hard margin has (\*) constrain but not (\*) (\*) so hard margin support vector =  $x_i$  where  $\alpha_i > 0 = x_i$  on margin

For Soft margin  $\rightarrow$  (\*)  $\rightarrow \alpha_i(y_i(w^T x_i + b) - 1 + \varepsilon_i) = 0 \rightarrow y_i(w^T x_i + b) - 1 + \varepsilon_i = 0 \rightarrow y_i(w^T x_i + b) = 1 - \varepsilon_i = y_i(w^T x_i + b) = 1 = \text{true}$  if  $\varepsilon_i = 0$

$\therefore$  Support Vector for SM is when  $\alpha > 0 \rightarrow$  contribute to weight (support vector)

If  $\alpha < C$ , then  $\varepsilon_i$  must be 0 and such point lies on the margin

so support vector is point

where  $0 < \alpha \leq C$

If  $\alpha = C$ , then  $\varepsilon_i > 0$  and point lies inside the margin and/or wrong side. Note  $\varepsilon_i = 0$  on margin/correct side  $\varepsilon_i > 0$  inside margin/wrong side. Include all  $x$  on and inside margin.

Support Vector where  $0 < \alpha < C$  has  $\varepsilon = 0$ , hence

$$y_i \left( \sum_{j|x_j \in S} \alpha_j y_j k(x_i, x_j) + b \right) = 1 \text{ and bias can be recovered } b = y_i - \sum_{j|x_j \in S} \alpha_j y_j k(x_i, x_j)$$

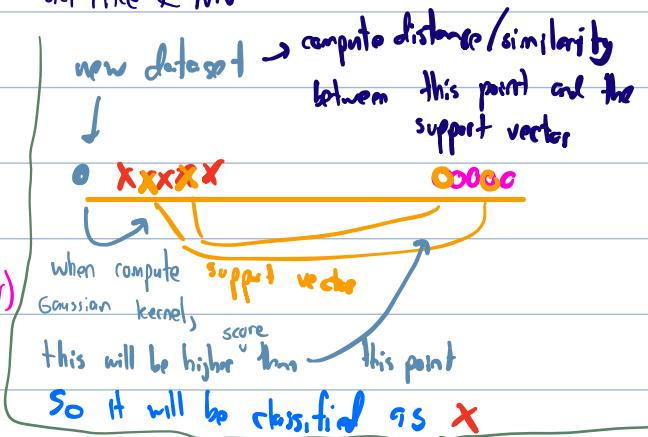
and new point  $x$  can be classified as  $h(x) = \text{sign} \left( \sum_j \alpha_j y_j k(x_j, x_k) + b \right)$  act like K-NN

## Multiple Class SVM

One-vs-rest : Train  $C$  SVM model for  $C$  classes, for each SVM of class  $x$ , predict if data belong to class  $x$  or other (all the rest of the class). (Predict either class  $x$  or other) the winner is then the class, where the distance from the hyperplane is maximal

One-vs-one : Train  $\binom{C}{2}$  classifiers (all possible pairing) and evaluate all.

Winner is the class with the majority vote (votes are weighted according to the distance from the margin)



$C_1$	$C_2$	Result
1	2	1
2	3	3
1	3	1

$C=2 = \text{win}$

## Lecture 10-11 - Dimensional Reduction

### Unsupervised Learning - modeling $p(x)$

generative transformation

- find the hidden (latent) structure in the data (latent distri.  $p(z)$  and  $p(x|z)$ , latent  $z$  usually unknown, has to be estimated)
- e.g. Clustering, Anomaly detection
- Can be viewed as Compression (higher to lower dim)

### Dimensionality Reduction

Idea : Try to reduce dimensionality while avoid info loss

#### Feature Selection

Choose "good" dim using prior knowledge or appropriate heuristic (remove low variance dim)

Dim Reduction via Linear Transf. → Change of basis (orthogonal basis trans.) and discard low variance dim.

Transform w/ orthogonal basis → preserve length  $F^T F = F F^T = I$

$X' = X \cdot F$  is the matrix containing all the transform point

Mean vector  $\bar{x}' = \bar{X} \cdot F \rightarrow \bar{X} \cdot F = \frac{1}{N} \sum_i x_i \cdot F = \frac{1}{N} \sum_i (x_i \cdot F) = \frac{1}{N} \sum_i x'_i = \bar{x}'$

$$\Sigma_x = F \cdot \Sigma_{\tilde{x}} \cdot F^T \rightarrow P^T \cdot (\frac{1}{N} (X^T X - \bar{X} \bar{X}^T)) \cdot F = \frac{1}{N} \Sigma (X_F)^T (X_F) = \frac{1}{N} \Sigma X^T X' = \Sigma_x$$

○ has zero mean

## PCA

Goal : Transform data s.t.  $\text{Cov}(x_i, x_j) = 0$  for all  $i \neq j$  (so data points are not correlated anymore)

Approach 1. Center the data ( $\tilde{x} = x - \bar{x}$  where  $\bar{x}$  = mean of each column  $\tilde{x}$  = zero mean)

2. Compute  $\Sigma_{\tilde{x}} = \begin{bmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \dots & \text{Cov}(x_1, x_d) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) & & \\ & \ddots & \ddots & \\ & & & \text{Var}(x_d) \end{bmatrix} = \frac{1}{N} X^T X$

(linearly independent)

$$\text{Var}(x_j) = \frac{1}{N} \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2$$

$$\text{Cov}(x_{j1}, x_{j2}) = \frac{1}{N} \sum_{i=1}^N (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)$$

3. Compute Eigenvector Decomposition of  $\Sigma_{\tilde{x}}$  to transform the coordinate

$$\Sigma_{\tilde{x}} = \Gamma \Lambda \Gamma^T \text{ where } \Gamma = \text{eigenvector matrix, } \Lambda = \text{diagonal matrix w/ eigenvalue } \lambda_i$$

Transform data =  $Y = \tilde{X} \cdot \Gamma$  and  $\Sigma_Y = \Gamma^T \Sigma_{\tilde{x}} \Gamma = \Gamma^T (\Gamma \Lambda \Gamma^T) \Gamma = \Lambda$  is the new covariance matrix of new coordinate system

$$\text{Variance of new system} = \Sigma_Y = \Lambda = \lambda_i \text{ and } \text{Cov}(Y_i, Y_j) = 0 \quad \forall i \neq j$$

Coordinate w/ low variance (low  $\lambda_i$ ) can be ignored  $\rightarrow$  only keep  $\Gamma$  col correspondingly to largest  $K$ -th  $\lambda_K$

$$\text{Pick } K \text{ s.t. } 90\% \text{ of var is explain} \rightarrow \sum_{i=1}^K \lambda_i \geq 0.9 \sum_{i=1}^d \lambda_i$$

## Alternative View of PCA

1. Maximize Variance
2. Minimize Residual
3. Projected data

Low rank approx - approx. original data A by a low rank matrix B

$$B \text{ has lower rank when approx. } A = \begin{bmatrix} 1.01 & 2.05 & 0.9 \\ -2.1 & -3.05 & 1.1 \\ 2.99 & 5.01 & 0.3 \end{bmatrix} \approx \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix} = B \quad \text{rank}(A)=3 \quad \text{rank}(B)=2$$

So by approx., when we rewrite coordinate in new system, it require less coordinate to describe. (b/c less dim)

$$\text{Best} = \min_{\text{error}, \text{rank}(B)=K} \|A - B\|_F^2 = \sum_{i=1}^n \sum_{j=1}^d (a_{ij} - b_{ij})^2 \text{ (error btw original & new approx. matrix)}$$

## SVD

$$A = U \Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T \quad A \in \mathbb{R}^{n \times d} \quad U \in \mathbb{R}^{n \times r} \quad \Sigma \in \mathbb{R}^{r \times r} \quad V \in \mathbb{R}^{d \times r} \quad r = \min(n, d)$$

Reading the  $U, \Sigma, V$  significant

$U = \text{span}(C(A)) \rightarrow$  column of A  $\rightarrow$  how each data row (eg. user) effect on concept (r col)

$V = \text{span}(C(A^T)) \rightarrow$  row of A  $\rightarrow$  how each feature (eg. movie) effect on concept (r col)

$\Sigma = \text{strength of each concept}$

$P = U \Sigma = AV$  give coordinate of points after being projected on the projection axis (new coordinate system)

$V$  act as a transformation matrix (like  $\Gamma$  in PCA)

Approx. matrix by setting smallest  $\sigma_{k+1} \dots \sigma_r$  to 0 and get best k approx.

→ choose k s.t. we can still explain var  $\geq 90\%$ .  $\sum_{i=1}^k \sigma_i^2 \geq 0.9 \sum_{i=1}^r \sigma_i^2$

Some Math Fact:  $\|x\|_F = \|x^T\|_F$   $\|x\|_F^2 = \text{trace}(x^T x)$   $\|ux\|_F^2 = \|vx\|_F^2 = \|x\|_F^2$  if  $U \approx V$ : orthogonal matrix

$$\|A\|_F^2 = \text{tr}((V\Sigma U^T)(U\Sigma V^T)) = \text{tr}(V\Sigma V^T) \cdot \text{tr}(V\Sigma^2 V^T) = \text{tr}(\Sigma^2) = \|\Sigma\|_F^2 = \sum_{i=1}^r \sigma_i^2$$

$$\text{if rank}(B)=k \quad A-B = \sum_{i=k+1}^r \sigma_i u_i v_i^T \quad \text{and} \quad \|A-B\|_F^2 = \sum_{i=k+1}^r \sigma_i^2$$

$$\text{In PCA, } \Sigma_x = x^T x = (U\Sigma V^T)^T (U\Sigma V^T) = V\Sigma^2 V^T = \Gamma \Lambda \Gamma^T \quad \text{if } V=\Gamma \text{ and } \Sigma^2 = \Lambda$$

can use eigenvalue decomposition to calculate the singular decomposition

## Matrix Factorization

System like Recommendation System has a lot of missing data

if we use SVD, we have to fill those missing value → e.g. fill w/ 0 → might be misleading

approx. SVD need to sum over all entry → What if we sum over existing entries

$$P \approx Q P^T \quad \min_{P, Q} \sum_{(u, i) \in S} (r_{ui} - q_u p_i^T)^2 \quad q_u = \text{row of } Q \quad p = \text{row of } P \quad (Q, P \text{ doesn't need to be orthogonal or unit length anymore})$$
$$Q = U\Sigma \quad P = V$$

## Alternating Optimization

• Pick initial  $P^{(0)}$  and  $Q^{(0)}$ ,  $t=0 \rightarrow$  missing entries are replace w/ 0 (or mean)

• Alternating keep one variable fix and optimizing for the other

$$P^{(t+1)} = \underset{P}{\arg \min} f(P, Q^{(t)}) \quad \text{fix } Q$$
$$f(P, Q) = \sum (r_{ui} - q_u p_i^T)^2$$
$$Q^{(t+1)} = \underset{Q}{\arg \min} f(P^{(t)}, Q)$$

• Repeat until converge (minimize sum of square Error  $f(P, Q)$ )

$$P^{(t+1)} = \sum_i \underset{p_i \in S, j \neq i}{\arg \min} \sum_{(u, j) \in S} (r_{uj} - q_u p_i^T)^2 = \sum_{(u, 1) \in S} (r_{u1} - q_u p_1^T)^2 + \sum_{(u, 2) \in S} (r_{u2} - q_u p_2^T)^2 + \dots = \text{is an OLS regression prob.}$$

only user u who have entry at col i

$$\min_{w \in \mathbb{R}} \sum_{i=1}^n (y_i - w^T x_i)$$

## Use SGD on Matrix Factorization instead

$$L = \sum_{(u, i) \in S} (r_{ui} - q_u p_i^T)^2$$

• Pick random u and i w/ rating  $r_{ui}$  (batch = 1)

• Compute gradient w.r.t. the parameter  $\frac{\partial L}{\partial q_u}, \frac{\partial L}{\partial p_i}$

• Update parameter  $q_u \leftarrow q_u - \eta \frac{\partial L}{\partial q_u}, p_i \leftarrow p_i - \eta \frac{\partial L}{\partial p_i}$

## Predict Missing value

Now that we have matrix  $Q, P$ , we can do matrix mult. to find missing rating in matrix  $R$

Cold start problem → new user has no data → Ask user to choose movie they like to create data

## Regularization

Overfit - too many parameters, not enough data points (too many missing data)

$$\min \sum_{(u,i) \in S} (r_{ui} - q_u p_i^T)^2 + \underbrace{\lambda_1 \|q_u\|^2}_{\text{"error"}} + \underbrace{\lambda_2 \|p_i\|^2}_{\text{"length"}}$$

If matrix has more data, the error term will dominate, so min problem will try to optimize error term more

If matrix has less data, the regularization will have more effect and it will try to reduce regularization term more

for each iteration in Alternative Opt.  $p^{(t+1)} = \arg \min_p \sum_i (r_{ui} - q_u p_i^T)^2 + \lambda_2 \|p_i\|^2$

for each iteration in SGD  $q_u \leftarrow q_u - \eta \left( \frac{\partial L}{\partial q_u} + \frac{\lambda}{2} \|q_u\|^2 \right) = q_u - \eta \left( \frac{\partial L}{\partial q_u} + 2\lambda_1 q_u \right)$

L1 - enforce sparsity → better interpretation (only few values contribute to result), less storage, faster processing

## Non-Negative Matrix Factorization

SVD might lead to factors containing negative values → hard to interpret

Constrained Optimization -  $\min_{P \geq 0, Q \geq 0} \|A - QP^T\|_F$

## Neighbor Graph Methods

PCA tries to find global structure

- Far away point can become nearest neighbor
- Loses clustering structure (local similarity) after projected

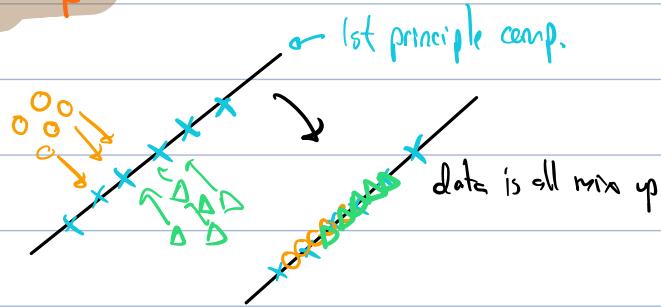
Idea: Preserve local structure - low dimensional neighborhood similar to the original neighborhood

1. Construct neighbor graph of high dim data

2. Initialize points in low dim space → construct neighborhood graph of this low dim data

3. Optimize coordinates in low dim space s.t. the neighborhood graph is similar

(move point around iteratively s.t. the points are group together again like in high dim data)



tsNE - + Distributed Stochastic Neighbor Embedding  $\rightarrow$  perform non linear dimensionality reduction

Minimize the KL-Divergence (measure of dissimilarity)  $\rightarrow$  use gaussian kernel

$$\min_{y_i} \text{KL}(P\parallel Q) = \sum_i \sum_{j \neq i} p_{ij} \log \left( \frac{p_{ij}}{q_{ij}} \right)$$

$y_i$  move  $y_i$  around

0 if distribution is identical  $p_{ij} = q_{ij}$

where  $p_{ij}$  = high dim similarity for input  $x_i$  and  $x_j$

$q_{ij}$  = low dim similarity for parameter  $y_i$  (to be optimized) and  $y_j$

We try to make distribution of higher & lower dim neighbor graph as similar as possible

Perplexity -  $\sigma^2$  - Variance of distribution to be consider = Area of neighborhood to preserve

## Autoencoder

PCA/SVD can only capture linearity

Autoencoder is a NN that find compact representation of data by learning to reconstruct its input

Goal : Minimize reconstruction error  $\min_w \frac{1}{N} \sum_{i=1}^N (f(x_i, w) - x_i)^2$

Find latent representation  $z \in \mathbb{R}^L$   $\dim(z) \ll \dim(x) \rightarrow$  perform non linear dimensionality reduction

$f_{enc}(x) = z$  encoder = project data to lower dim

$f_{dec}(z) \approx x$  decoder = reconstruct data from the latent code

$L \ll D$  allow autoencoder to capture the important feature of the data

$L \geq D$  make autoencoder learn just an identity function (learn to copy every pixel)

## Lecture 12: Clustering

Dimensionality Reduction - project high dim data into low dim space such that similarity stay the same ( $z \in \mathbb{R}^k$ )

(Clustering - Group object  $x_i$  into  $K$  cluster based on their similarity ( $z \in$  discrete number))

Goal : Assign  $x_i$  to one of the  $K$  clusters. s.t. (in case of non overlapping cluster)

- Objects within same group are similar to each other

- Object btw diff group are dissimilar from each other

## K-Means Clustering (based on distance)

Similarity = use distance function

1. Manhattan distance  $\|x_i - x_j\|_1$

2. Euclidean distance  $\|x_i - x_j\|_2$

Mahalanobis  $\sqrt{(x_i - \bar{x})^T \Sigma^{-1} (x_i - \bar{x})}$  generalisation of  $L_2$  dist  
if  $\Sigma^{-1} = I$  then we get  $L_2$

$\mu_k$  = centroid of cluster  $k$        $z_{ik}$  = cluster indicator =  $x_i$  belong to cluster  $k$

K-means objective : Find best centroid  $\mu$  and cluster assignments  $Z$

$$\min J(X, Z, \mu) = \min_{Z, \mu} \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|x_i - \mu_k\|_2^2$$

Use Alternating Optimization to update  $Z$  and  $\mu$  in turn

1. initialise centroid  $\mu$

2. Update cluster indicator  $z_{ik}$  (solve  $\min_z J(X, Z, \mu)$ )

calculate distance btw  $x_i$  and  $\mu_k$ . Assign  $x_i$  to cluster with least distance

3. Update Centroid (solve  $\min_{\mu} J(X, Z, \mu)$ )

$\mu_k = \frac{1}{N_k} \sum_{i=1}^N z_{ik} x_i$      $N_k$  = # of data in cluster    compute mean only count data that belong to cluster.

Iterate btw 2&3

- initial location of  $\mu$  effect perf.  $\rightarrow$  if far away from data can lead to bad result

K-Means ++ Algorithm  $\hookrightarrow$  solution

1. Choose  $\mu_1$  at random among the data points

2. For each point  $x_i$  compute  $D_i^2 = \|x_i - \mu_1\|_2^2$

3. Sample next  $\mu_2$  from  $\{x_i\}$  w/ probability proportional to  $D_i^2$

so the further away  $x_i$ , the more likely it's to get chosen as the next  $\mu$

4. Recompute  $D_i$  for all centroid  $\mu$      $D_i^2 = \min \{\|x_i - \mu_1\|_2^2, \|x_i - \mu_2\|_2^2\}$

5. Repeat 3&4 to find next  $\mu$  until get all  $K$  centroid.

Limitation of K-Mean

can't detect cluster w/ overlapping convex hull, sensitive to outliers, no uncertainty measure, sensitive to initialisation.

Gaussian Mixture Model (GMM) (Probabilistic model)

Goal : Design probabilistic model for data  $p(x|\theta)$ , estimate  $\theta$

Latent Variable marginalization  $\rightarrow p(x) = \sum_j p(x, z_j)$

$$p(x|\theta) = \sum_z p(x, z|\theta) = \sum_z p(x|z, \theta) p(z|\theta) \text{ and then } \theta^* = \arg \min_{\theta} p(x|\theta)$$

$z$  are never observe and only use to simplify model = latent variable

$$p(x, z|\theta) = \underbrace{p(x|z, \theta)}_{\text{cluster}} \underbrace{p(z|\theta)}_{\text{likelihood prior}}$$

$$\text{where } p(z|\theta) = \text{category}(\pi) \quad \begin{matrix} \text{prob for each cluster} \\ \pi = [0.2 \ 0.5 \ 0.3] \end{matrix}$$

$$\text{Param } \theta = \{\pi, \mu, \Sigma\}$$

given  $x$  belong to

cluster  $k$

param for each cluster  $k$

$$p(x|z_k=1, \theta) = N(x | \mu_k, \Sigma_k)$$

## Generative Process of GMM

We have  $K$  cluster, known parameter  $\{\mu_k, \Sigma_k\}$  and  $\pi$  for all cluster

1. Sample cluster  $z_k$  based on  $\pi \rightarrow z \sim \text{Cat}(\pi)$

2. Get  $\mu_k, \Sigma_k$  from  $z_k$  then we can sample  $x$  from the likelihood distribution  $x \sim N(\mu_k, \Sigma_k)$  for  $z_k$

What if we know  $x$  but not  $\pi, \mu, \Sigma$

$$p(x|\pi, \mu, \Sigma) = \sum_{k=1}^K p(x, z_k=1 | \pi, \mu_k, \Sigma_k) = \sum_{k=1}^K p(z_k=1 | \pi) p(x | z_k=1, \mu_k, \Sigma_k) = \sum_{k=1}^K \pi_k \cdot N(x | \mu_k, \Sigma_k)$$

GMM Likelihood  $= p(x|\pi, \mu, \Sigma) = \sum_{k=1}^K \pi_k \cdot N(x | \mu_k, \Sigma_k)$  = weighted sum of gaussian distribution

$$\text{Log likelihood } \log p(x|\pi, \mu, \Sigma) = \sum_{i=1}^N \log \left( \sum_{k=1}^K \pi_k \cdot N(x_i | \mu_k, \Sigma_k) \right)$$

Learning -  $\pi^*, \mu^*, \Sigma^* = \arg \max_{\pi, \mu, \Sigma} \log p(x|\pi, \mu, \Sigma)$

Inference - Given  $\{\pi, \mu, \Sigma\}$ , we want to assign point to cluster

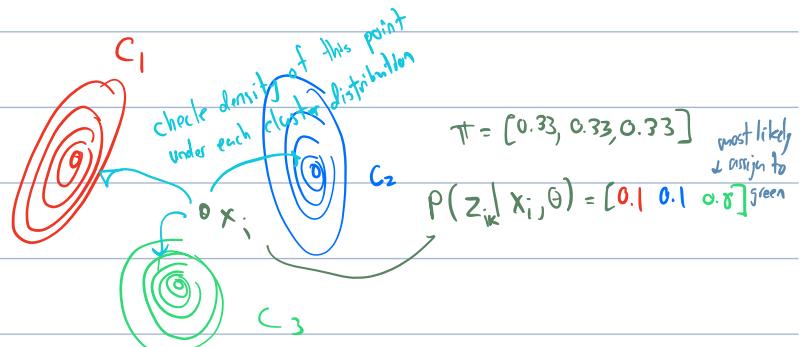
- Compute Posterior Distribution of cluster indicator  $z$   $P(z|x, \pi, \mu, \Sigma)$

## Inference in GMM

$$p(z_{ik}=1 | x_i, \pi, \mu, \Sigma) = \frac{p(z_{ik}=1 | \pi) p(x_i | z_{ik}=1, \mu_k, \Sigma_k)}{p(x_i | \pi, \mu, \Sigma)} = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_i | \mu_j, \Sigma_j)}$$

We call this "responsibility" of component  $k$  for observation  $i$

$\gamma(z_{ik}) = p(z_{ik}=1 | x_i, \pi, \mu, \Sigma) \rightarrow$  prob of assigning to the cluster given data



# Learning in GMM

add  $z$  to the  $p(x|\theta)$  to make it easier to solve

$$\pi^*, \mu^*, \xi^* = \max_{\pi, \mu, \xi} \log p(x, z | \pi, \mu, \xi) \Rightarrow \text{Estimate } \pi^{(t+1)}, \mu^{(t+1)}, \xi^{(t+1)} = \max_{\pi, \mu, \xi} \mathbb{E}_{z \sim \gamma_t(z)} [\log p(x, z | \pi, \mu, \xi)]$$

## Use EM Algorithm (Expectation Maximization)

### E Step

Compute & Update responsibility  $\gamma$  (posterior distribution) based on initial/previous  $\pi_+, \xi_+, \mu_+$

$$\gamma_+(z_{ik}) = p(z_k | x_i, \pi_+, \mu_k, \xi_k) = \frac{p(z_k | \pi_+) p(x_i | \mu_k, \xi_k)}{p(x_i | \pi_+, \mu_k, \xi_k)} = \frac{\pi_k N(x_i | \mu_k, \xi_k)}{\sum_{j=1}^K \pi_j^{(t)} N(x_i | \mu_j, \xi_j)} \rightarrow \text{get } \gamma_+(z_{ik}) \text{ to be use in M step}$$

### M-Step - update $\theta$

Maximize Expectation  $\theta^* = \max_{\theta} \mathbb{E}(\log p(x, z | \theta)) \rightarrow \text{get optimal } \pi^*, \mu^*, \xi^*$

$$\theta^{(t+1)} = \max_{\theta} \mathbb{E}_{z \sim \gamma_t(z)} [\log p(x, z | \theta)] = \max_{\theta} \sum_{k=1}^K \sum_{i=1}^N \gamma_+(z_{ik}) \log (x_i, z_{ik} | \theta) \quad \text{use } \gamma_+(z) \text{ that was calculate from previous step}$$

$$M_K^{(t+1)} = \frac{1}{N_K} \sum_{i=1}^N \gamma_+(z_{ik}) x_i \quad \xi_K^{(t+1)} = \frac{1}{N_K} \sum_{i=1}^N \gamma_+(z_{ik}) (x_i - M_K^{(t+1)})^T \quad \pi_K^{(t+1)} = \frac{N_K}{N}$$

get  $\mu^{(t+1)}, \xi^{(t+1)}, \pi^{(t+1)}$  → use in E step to compute new  $\gamma_{t+1}(z_{ik}) \rightarrow$  get  $\theta^{(t+2)}$  from M-step → iterate until converge

## Kinda Explanation of EM ... ? (go read pg. 30-34)

$$\text{Denote } L(q, \theta) = \mathbb{E}_{z \sim q} [\log \frac{p(x, z | \theta)}{q(z)}] = \mathbb{E}_{z \sim q} [\log p(x, z | \theta)] + H(q(\cdot)) \quad \text{indep. of } \theta$$

It's proved that  $L(q, \theta)$  is lower bound of  $\log p(x | \theta)$  for any  $q$  → any lower bound

If  $q(z) = \gamma(z) = p(z | x, \theta)$ , then we maximize  $L(q, \theta)$  and  $L(q, \theta) = \log p(x | \theta)$

From the maximum  $L(q, \theta)$  → we find  $\theta$  that maximize  $L(q, \theta) \rightarrow$  get  $\theta^{(t+1)}$  to use in  $L(q, \theta)$  and  $\gamma(z)$

to get the next tight lower bound

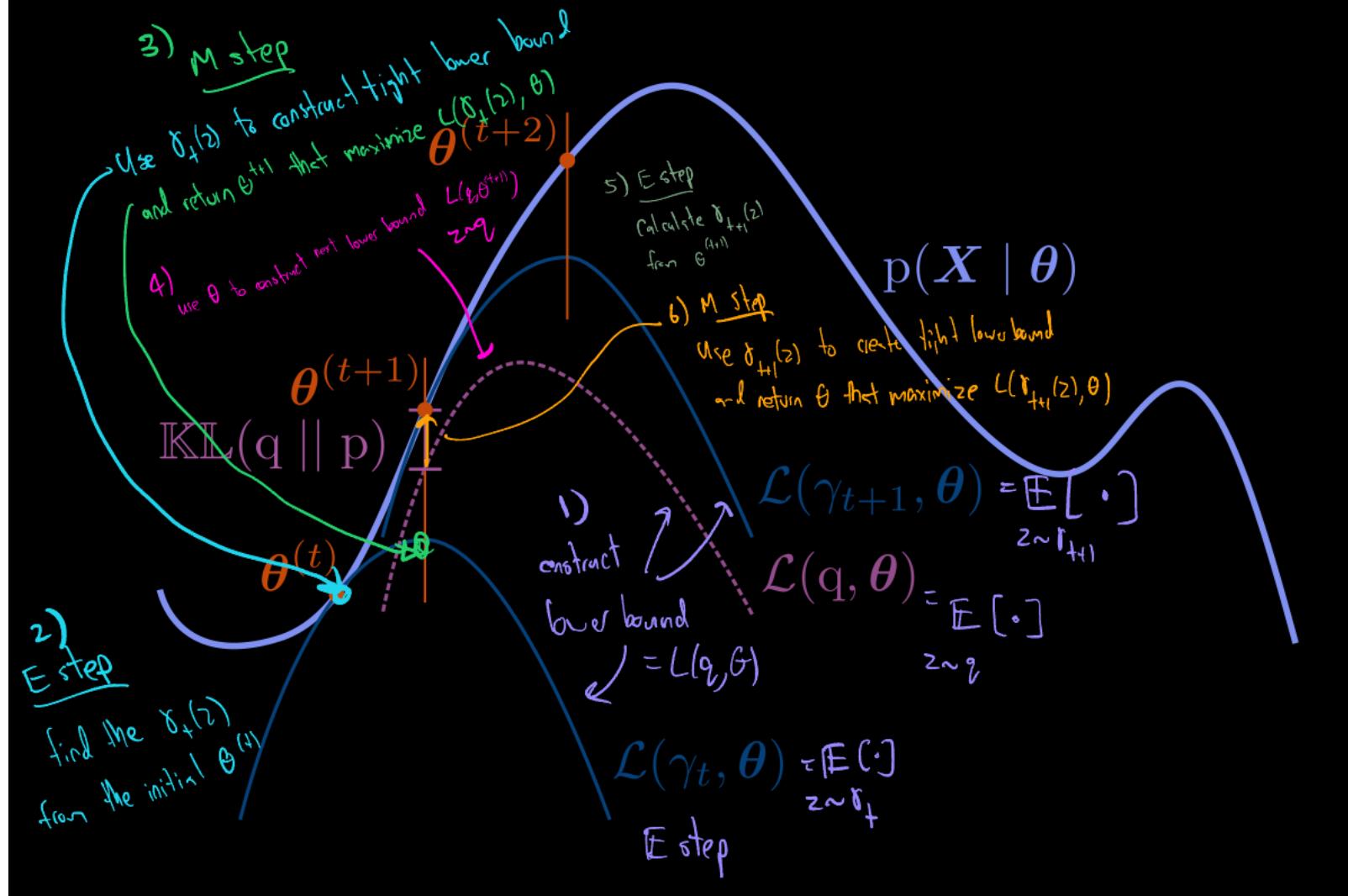
(tight lowerbound)

$\theta$  that have higher  $L(q, \theta)$

As a consequence -  $L(q, \theta)$  is push up against  $\log p(x | \theta)$  always

- If  $L(q, \theta)$  doesn't change, we are at local maximum of  $\log p(x | \theta)$  and EM converge

# EM algorithm: Illustration



## Hierarchical Clustering

Agglomerative - Bottom up, start with all cluster and start merging each pair until get 1 big cluster

Divisive - Top-down - start from 1 cluster and split it

How to choose pair (let  $x_i \in G, x_j \in H$  cluster G, H)

Complete-linkage  $\rightarrow \max d(x_i, x_j) \rightarrow$  look at max dist of sample b/w 2 cluster

Single Linkage  $\rightarrow \min d(x_i, x_j) \rightarrow$  min dist  $\xrightarrow{\text{---}}$

Centroid  $\rightarrow$  dist from center of cluster

Unweighted Average  $\rightarrow$  Mean sample distance  $\frac{1}{N_G N_H} \sum_{x_i \in G, x_j \in H} d(x_i, x_j)$

Weighted Average