

Machine Learning

NYU Shanghai
Spring 2020

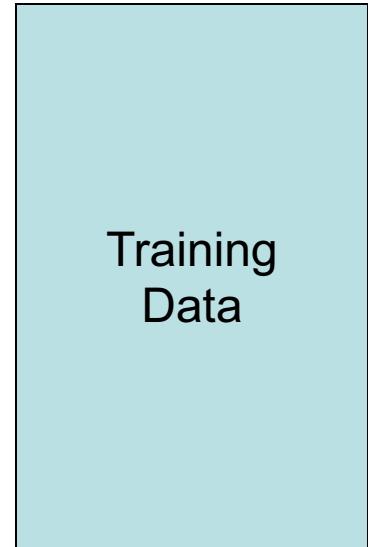
Some Popular Types of Machine Learning

- Supervised vs. unsupervised
 - Supervised often used in practice
 - But... unsupervised is closer to human learning
- Discriminative vs. generative
 - discriminative often used in practice
 - But... generative is more fun
- Deep Learning
 - Deep Learning is just a fancy term for Neural networks ++
- Reinforcement Learning

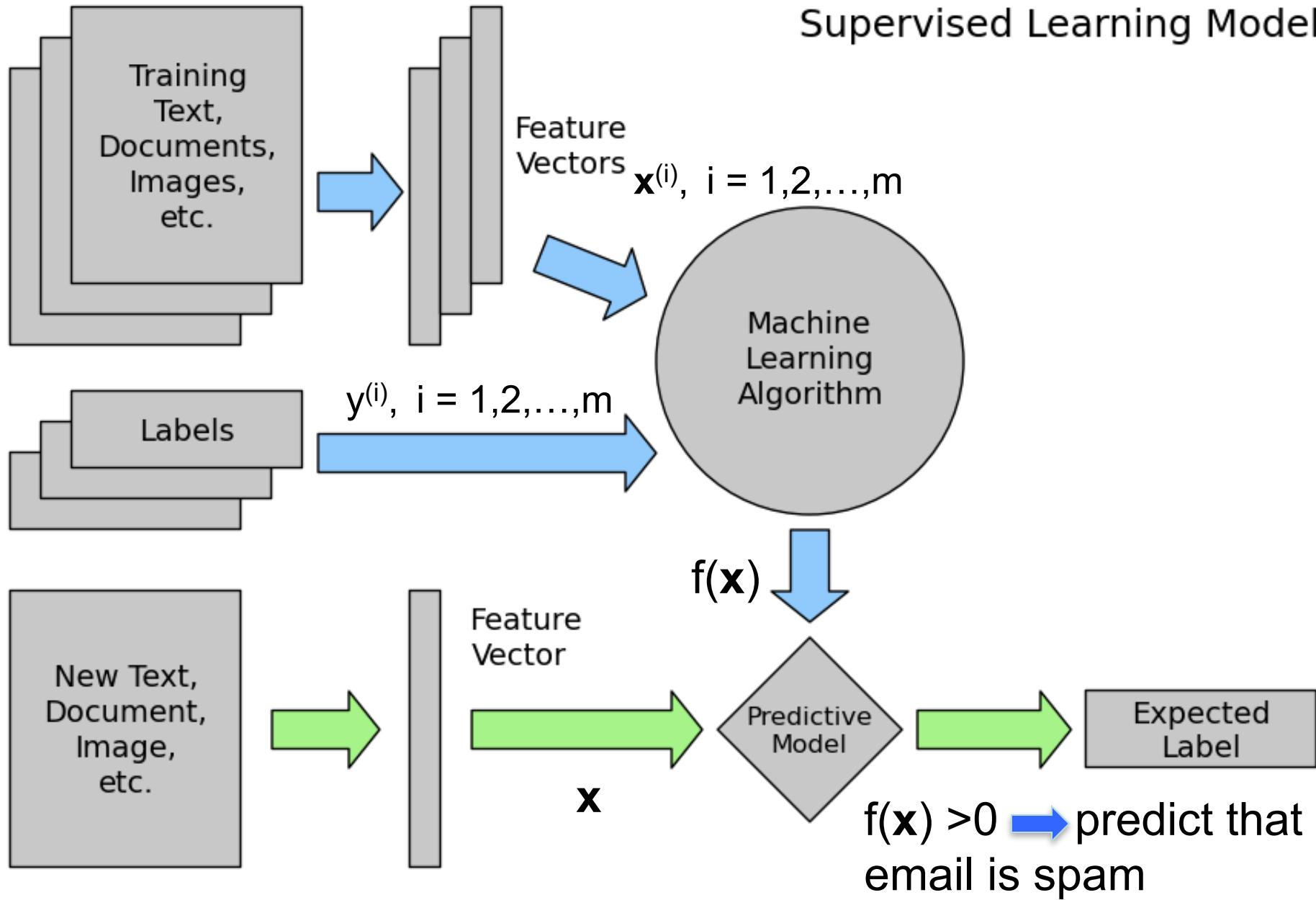
ML Methodology

- **Experimentation cycle**

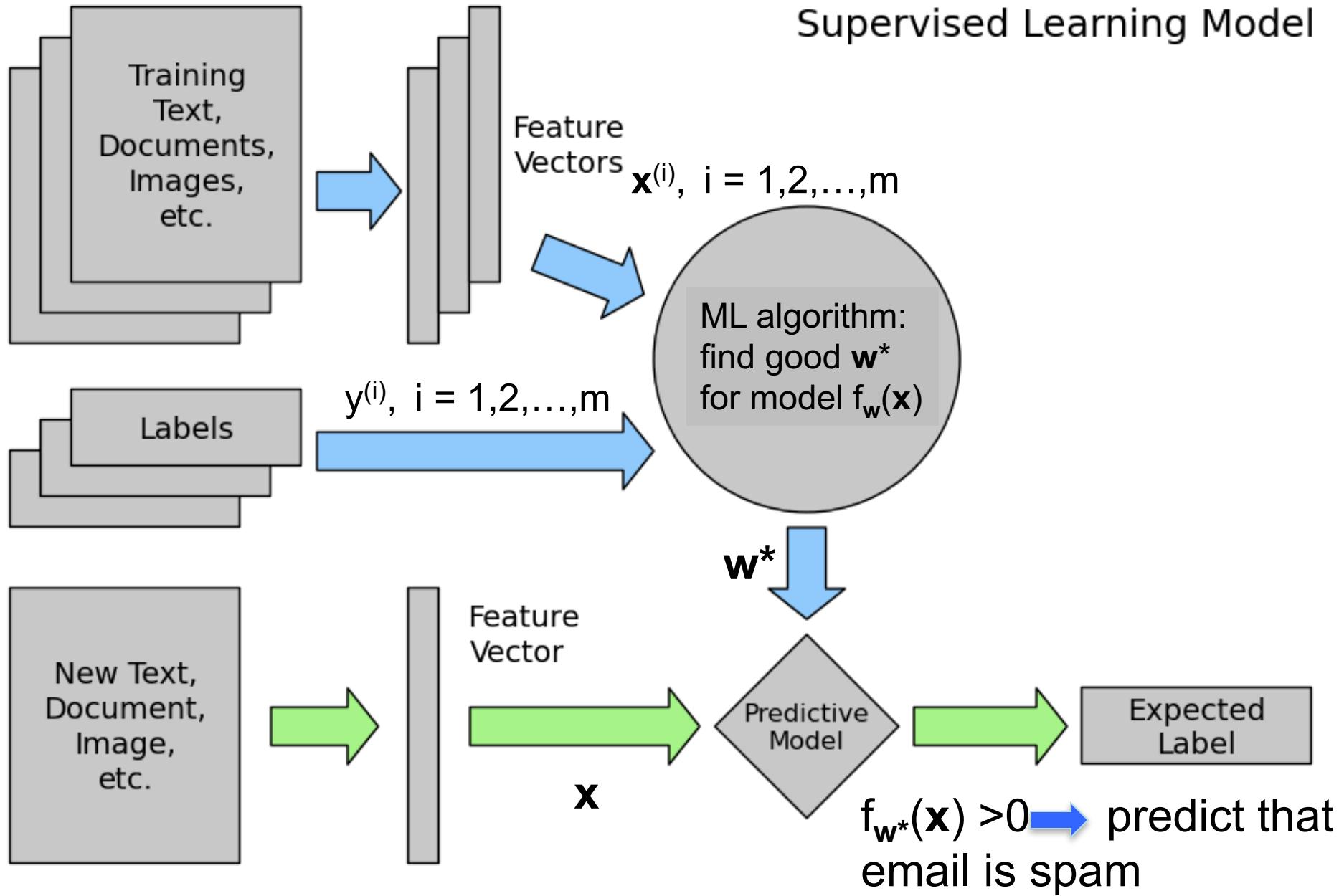
- Select model, features, algorithm
- Train using **training data** (ie, use algorithm to find good weights)
- Evaluate using **validation data**: number of instances in validation set predicted correctly.
- Modify features, hyperparameters
- Train again....
- Find best features, hyperparameters
- **Avoid Overfitting: to avoid a good performance on training data but a bad performance on validation data**
- Using bestfeatures and hyperparameters, train with both training data and validation data
- Compute final accuracy of test data
- Very important: never “peek” at the test set!



Supervised Learning Model



Supervised Learning Model



The Model

- Assumptions:
 - The model is a parameterized function $f_w(\mathbf{x})$, where w is a vector of parameters (also called weights).
 - $f_w(\mathbf{x})$ has a specific form. E.g., $f_w(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_nx_n$ is a *linear model*.
- Under the assumptions above, the learning process only determines what specific w we should use.
- We often use w^* to denote the best w

The perceptron algorithm

- Start with weight vector = $\mathbf{w} = (0, 0, \dots, 0)$
- Cycle through the training examples $\mathbf{x}^{(i)}, y^{(i)}$,
 $i=1,2,\dots,m,1,2,\dots,m,\dots$

- For $i=1,2,\dots,m,1,2,\dots,m,\dots$

– if $y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)}) > 0$: (example correctly classified)

$$\mathbf{w} = \mathbf{w}$$

– else: (example incorrectly classified)

$$\mathbf{w} = \mathbf{w} + y^{(i)} \mathbf{x}^{(i)}$$

Any problem?

这么强的假设，有没有副作用？

If for some \mathbf{w} , get $y^{(i)}(\mathbf{w} \cdot \mathbf{x}^{(i)}) > 0$ for all $i=1,\dots,m$, stop and use that \mathbf{w} !

what does all this have to do with
function approximation?

instead of $F: X \rightarrow Y$,
learn $P(Y | X)$

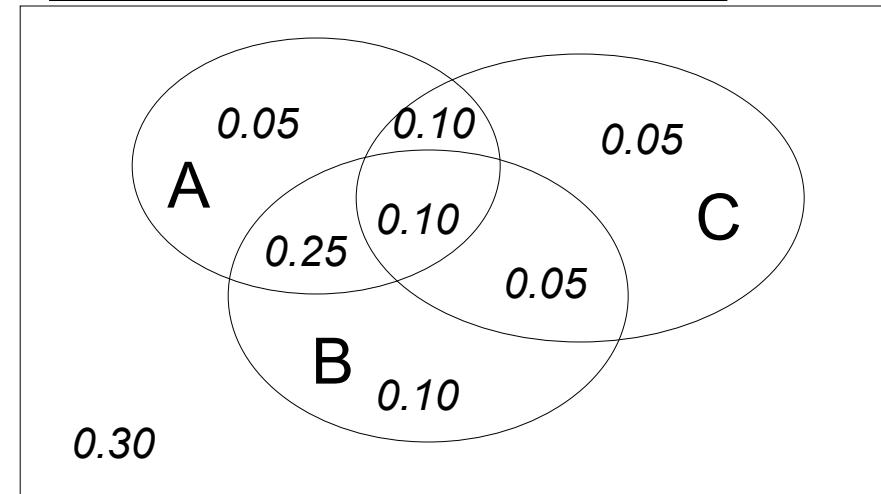
The Joint Distribution

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values (M Boolean variables $\rightarrow 2^M$ rows).
2. For each combination of values, say how probable it is.

Example: Boolean variables A, B, C

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10



[A. Moore]

sounds like the solution to
learning $F: X \rightarrow Y$,
or $P(Y | X)$.

Are we done?

Learning and the Joint Distribution

gender	hours_worked	wealth	
Female	v0:40.5-	poor	0.253122
		rich	0.0245895
	v1:40.5+	poor	0.0421768
		rich	0.0116293
Male	v0:40.5-	poor	0.331313
		rich	0.0971295
	v1:40.5+	poor	0.134106
		rich	0.105933

Suppose we want to learn the function $f: \langle G, H \rangle \rightarrow W$

Equivalently, $P(W | G, H)$

Solution: learn joint distribution from data, calculate $P(W | G, H)$

e.g., $P(W=\text{rich} | G = \text{female}, H = 40.5-) =$

$$\frac{P(W=r \wedge G=f \wedge H=40.5-)}{P(G=f \wedge H=40.5-)} = \frac{0.024}{0.277} \approx 0.09$$

[A. Moore]

sounds like the solution to
learning $F: X \rightarrow Y$,
or $P(Y | X)$. $2^{10} = 1024$

Main problem: learning $P(Y|X)$
can require more data than we have

consider learning Joint Dist. with 100 attributes
of rows in this table? $2^{100} \approx 10^{30}$
of people on earth? 10^9
fraction of rows with 0 training examples? 0.9999

Can we reduce params using Bayes Rule?

Suppose $X = \langle X_1, \dots, X_n \rangle$

where X_i and Y are boolean RV's

$$P(Y|X) = \frac{\cancel{P(X|Y)} P(Y)}{\cancel{P(X)}}$$

How many parameters to define $P(X_1, \dots, X_n | Y)$?

$$P(X | Y=1) \rightarrow 2^n - 1$$

$$P(X | Y=0) \rightarrow 2^n - 1$$

$$2^{'}$$

How many parameters to define $P(Y)$?

question: why don't we have to estimate $p(x)$?

Naïve Bayes

Naïve Bayes assumes

$$P(X_1 \dots X_n | Y) = \prod_i P(X_i | Y)$$

i.e., that X_i and X_j are conditionally independent given Y , for all $i \neq j$

Naïve Bayes uses assumption that the X_i are conditionally independent, given Y . E.g., $P(X_1|X_2, Y) = P(X_1|Y)$

Given this assumption, then:

$$\begin{aligned}P(X_1, X_2|Y) &= P(X_1|X_2, Y)P(X_2|Y) \\&= P(X_1|Y)P(X_2|Y)\end{aligned}$$

in general:
$$P(X_1 \dots X_n|Y) = \prod_i P(X_i|Y)$$

How many parameters to describe $P(X_1 \dots X_n|Y)$? $\underline{P(Y)}$?

- Without conditional indep assumption? $2(2^n - 1) + 1$
- With conditional indep assumption? $2n + 1$

Again, any problem?

Naïve Bayes: Subtlety #1

Often the X_i are not really conditionally independent

- We use Naïve Bayes in many cases anyway, and it often works pretty well
 - often the right classification, even when not the right probability (see [Domingos&Pazzani, 1996])
- What is effect on estimated $P(Y|X)$?
 - Extreme case: what if we add two copies: $X_i = X_k$

$$P(Y=1|X) = P(Y=1) P(X_1|Y=1) P(X_2|Y=1) \dots$$

Extreme case: what if we add two copies: $X_i = X_k$

we know it is harmful, what to do?

Naïve Bayes: Subtlety #2

If unlucky, our MLE estimate for $P(X_i | Y)$ might be zero.
(for example, $X_i = \text{birthdate}$. $X_i = \text{Jan_25_1992}$)

$$\theta = \hat{P}(X = 1 | S = 1)$$

- Why worry about just one parameter out of many?

$$P(Y=1 | X_1, \dots, X_n) = \frac{P(Y=1) \prod_i P(X_i | Y=1)}{P(X_1, \dots, X_n)}$$

- What can be done to address this?

→ MAP cst, mate!

Estimating Parameters: Y, X_i discrete-valued

Maximum likelihood estimates:

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\}}{|D|}$$

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_j | Y = y_k) = \frac{\#D\{X_i = x_j \wedge Y = y_k\}}{\#D\{Y = y_k\}}$$

MAP estimates (Beta, Dirichlet priors):

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\} + (\beta_k - 1)}{|D| + \sum_m (\beta_m - 1)}$$

Only difference:
“imaginary” examples

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_j | Y = y_k) = \frac{\#D\{X_i = x_j \wedge Y = y_k\} + (\beta_k - 1)}{\#D\{Y = y_k\} + \sum_m (\beta_m - 1)}$$

What if we have continuous X_i ?

Gaussian Naïve Bayes (GNB): assume

$$p(X_i = x | Y = y_k) = \frac{1}{\sqrt{2\pi\sigma_{ik}^2}} e^{-\frac{1}{2}\left(\frac{x - \mu_{ik}}{\sigma_{ik}}\right)^2}$$

Sometimes assume variance

- is independent of Y (i.e., σ_i),
- or independent of X_i (i.e., σ_k)
- or both (i.e., σ)

Gaussian Naïve Bayes Algorithm – continuous X_i (but still discrete Y)

- Train Naïve Bayes (examples)

for each value y_k

estimate* $\pi_k \equiv \underline{P(Y = y_k)}$

for each attribute X_i estimate $\underline{P(X_i|Y = y_k)}$

• class conditional mean μ_{ik} , variance σ_{ik}

- Classify (X^{new})

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new}|Y = y_k)$$

$$Y^{new} \leftarrow \arg \max_{y_k} \pi_k \prod_i \mathcal{N}(X_i^{new}; \mu_{ik}, \sigma_{ik})$$

* probabilities must sum to 1, so need estimate only n-1 parameters...

Estimating Parameters: Y discrete, X_i continuous

Maximum likelihood estimates:

$$\hat{\mu}_{ik} = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j X_i^j \delta(Y^j = y_k)$$

ith feature kth class jth training example
 $\delta()=1$ if $(Y^j=y_k)$
else 0

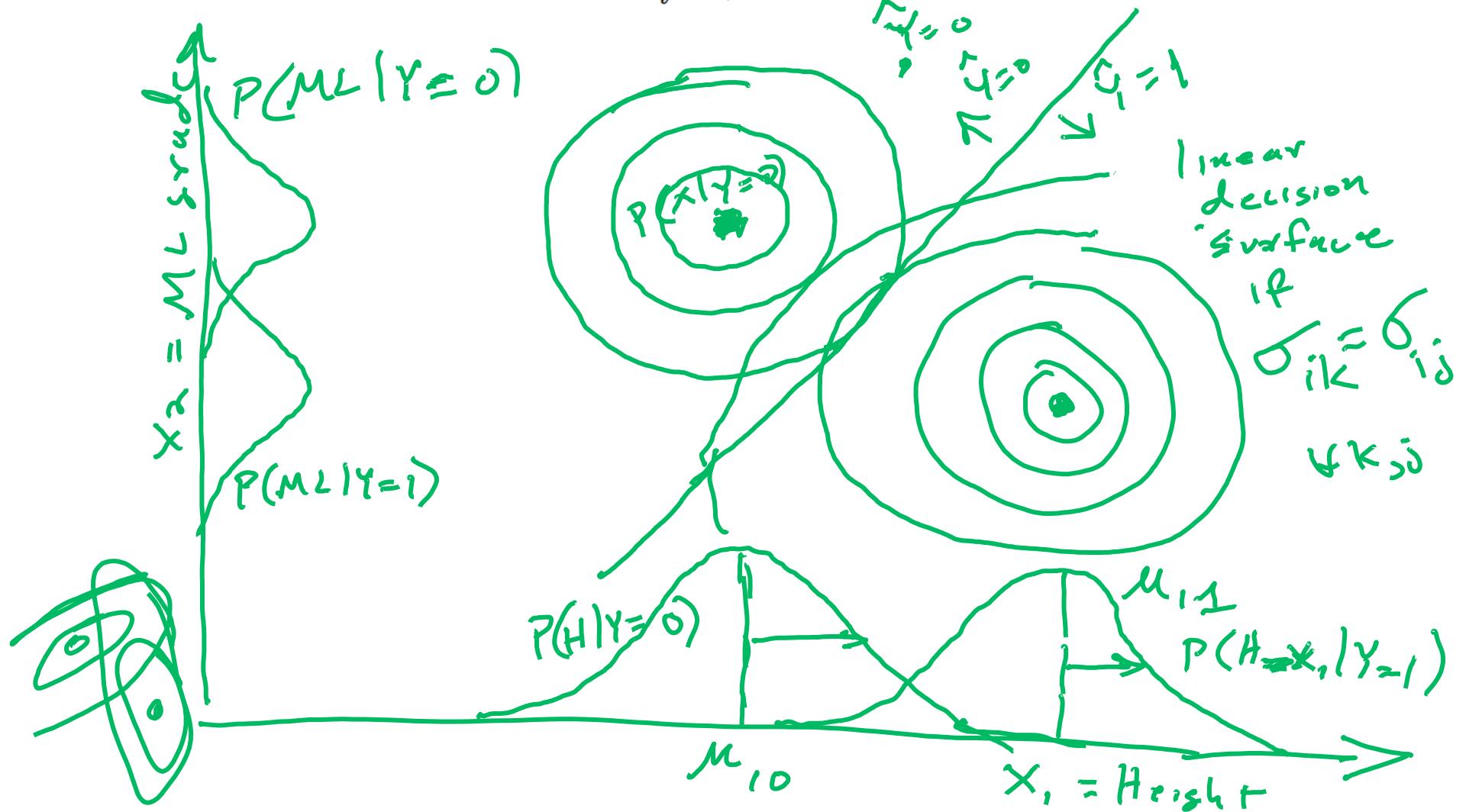
$$\hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j (X_i^j - \hat{\mu}_{ik})^2 \delta(Y^j = y_k)$$

we can have imaginary count here, too.

Gaussian Naïve Bayes – Big Picture

Example: $Y = \text{PlayBasketball}$ (boolean), $X_1 = \text{Height}$, $X_2 = \text{MLgrade}$

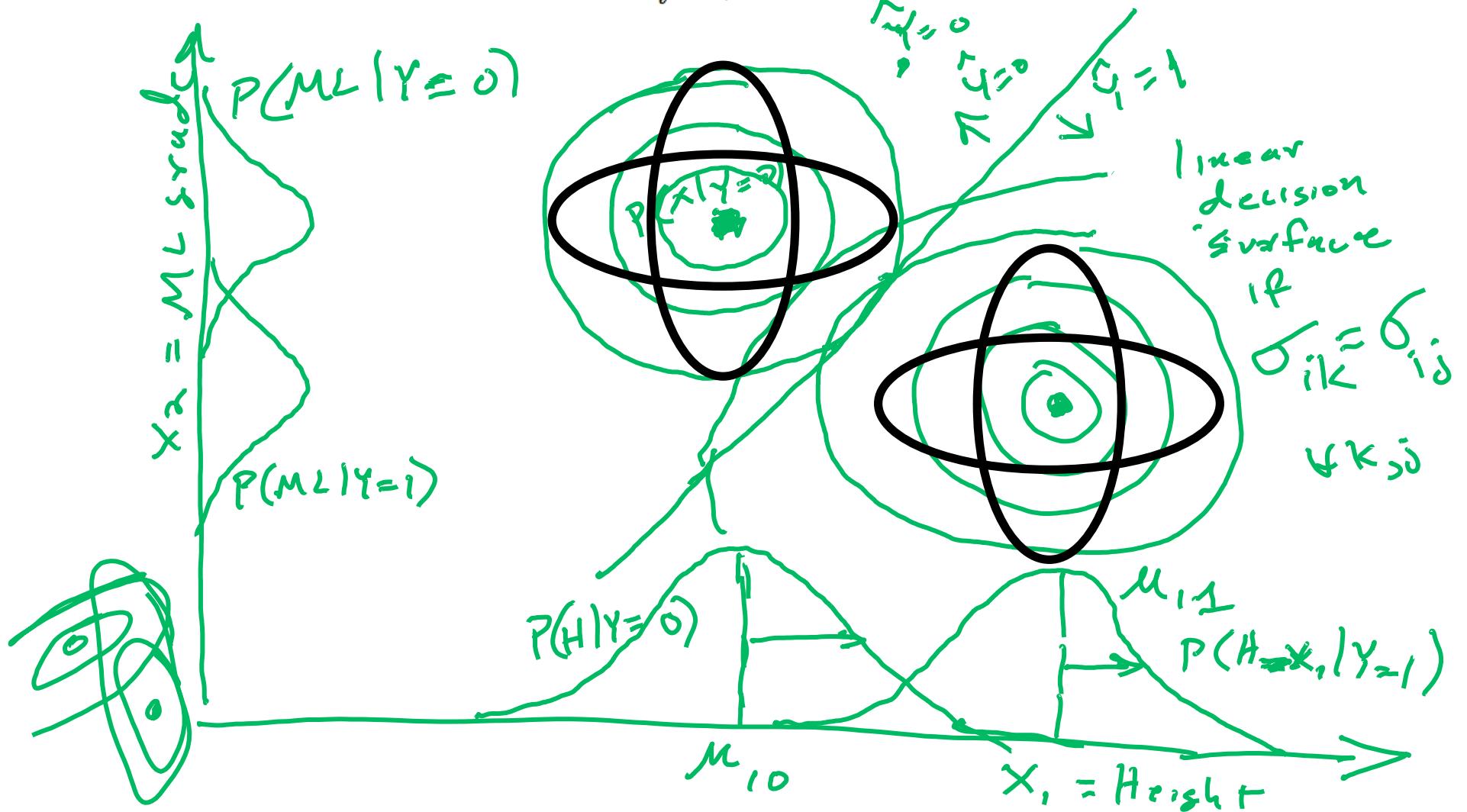
$$Y^{new} \leftarrow \arg \max_{y \in \{0,1\}} P(Y = y) \prod_i P(X_i^{new} | Y = y) \quad \text{assume } P(Y=1) = 0.5$$



Gaussian Naïve Bayes – Big Picture

Example: $Y = \text{PlayBasketball}$ (boolean), $X_1 = \text{Height}$, $X_2 = \text{MLgrade}$

$$Y^{new} \leftarrow \arg \max_{y \in \{0,1\}} P(Y = y) \prod_i P(X_i^{new} | Y = y) \quad \text{assume } P(Y=1) = 0.5$$



Logistic Regression

Idea:

- Naïve Bayes allows computing $P(Y|X)$ by learning $P(Y)$ and $P(X|Y)$
- Why not learn $P(Y|X)$ directly?

- Consider learning $f: X \rightarrow Y$, where
 - X is a vector of real-valued features, $\langle X_1 \dots X_n \rangle$
 - Y is boolean
 - assume all X_i are conditionally independent given Y
 - model $P(X_i | Y = y_k)$ as Gaussian $N(\mu_{ik}, \sigma_i)$ *not σ_{ik}*
 - model $P(Y)$ as Bernoulli (π)
- What does that imply about the form of $P(Y|X)$?

$$P(Y = 1 | X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

Very convenient!

$$P(Y = 1|X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

$$P(Y = 0|X = \langle X_1, \dots, X_n \rangle) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

$$\frac{P(Y = 0|X)}{P(Y = 1|X)} = \exp(w_0 + \sum_i w_i X_i)$$

linear
classification
rule!

implies

$$\ln \frac{P(Y = 0|X)}{P(Y = 1|X)} = w_0 + \sum_i w_i X_i$$

Back to Supervised Machine Learning

- $\mathbf{x} = (x_1, x_2, \dots, x_n)$ “feature vector”
- $f_w(\mathbf{x})$ is a model
- *For example, Linear Model:*
$$f_w(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$
where $\mathbf{w} = (w_0, w_1, w_2, \dots, w_n)$
- ML algorithm uses training set to find good w^* for $f_w(\mathbf{x})$
- Given new email \mathbf{x} , predict \mathbf{x} is SPAM if $f_{w^*}(\mathbf{x}) > 0$; otherwise not SPAM

Hyperplane

- Let $\mathbf{w} = (w_1, w_2, \dots, w_n)$ be any fixed vector and c any scalar
- The equation $w_1u_1 + w_2u_2 + \dots + w_nu_n = c$ (equivalently $\mathbf{w} \cdot \mathbf{u} = c$) defines a hyperplane in n -dimensional space.
- \mathbf{w} is perpendicular to plane
 - The hyperplane cuts \mathbb{R}^n into two halves
 - Points $\mathbf{w} \cdot \mathbf{u} > c$ lie on one side of plane
 - Points $\mathbf{w} \cdot \mathbf{u} < c$ lie on the other side

Derive form for $P(Y|X)$ for Gaussian $P(X_i|Y=y_k)$ assuming $\sigma_{ik} = \sigma_i$

$$P(Y=1|X) = \frac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)}$$

$$= \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}} \quad x = \exp(\ln x)$$

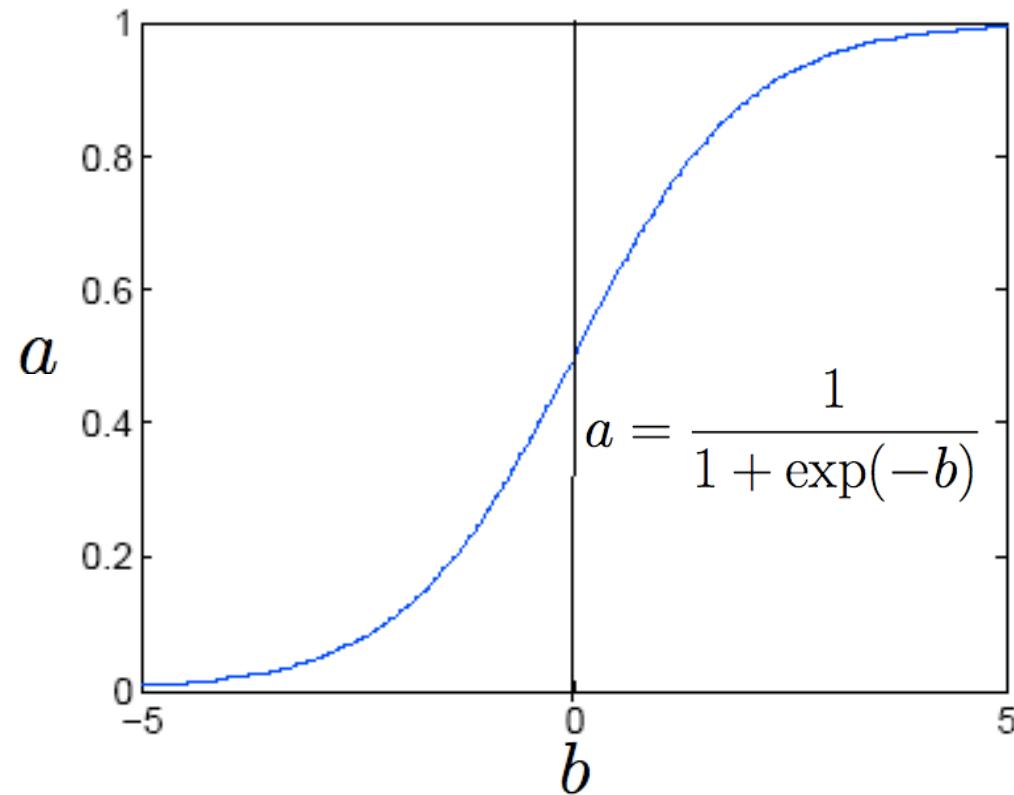
$$= \frac{1}{1 + \exp(\ln \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)})} \quad \pi \equiv P(Y=1)$$

$$= \frac{1}{1 + \exp(-\ln \frac{1-\pi}{\pi}) + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}} \quad c$$

$$P(x | y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} e^{\frac{-(x-\mu_{ik})^2}{2\sigma_{ik}^2}} \quad \sum_i \left(\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} X_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right)$$

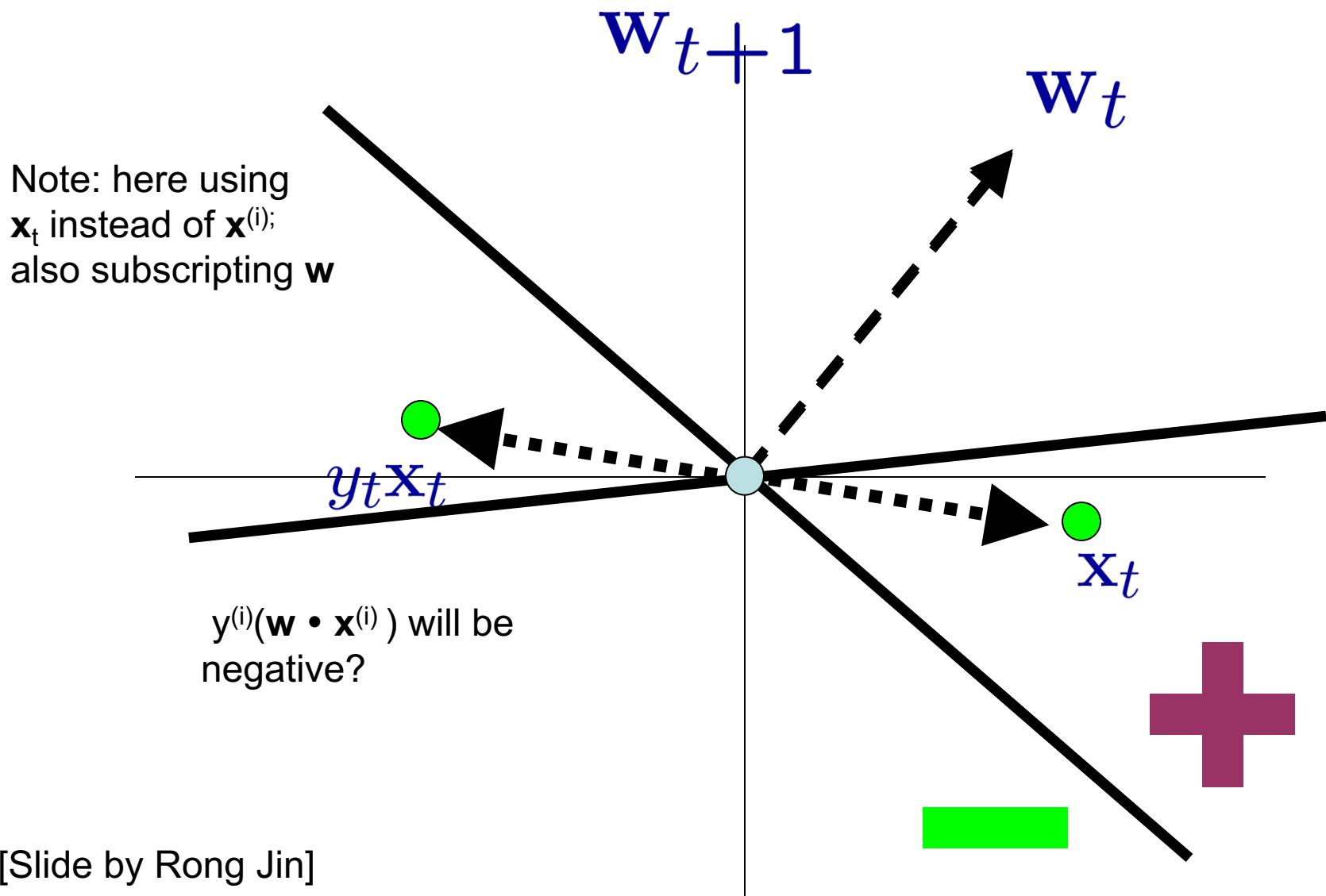
$$P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

Logistic function



$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

Geometrical Interpretation



Training Logistic Regression: MCLE

- Choose parameters $W = \langle w_0, \dots, w_n \rangle$ to maximize conditional likelihood of training data

where

$$P(Y = 0|X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1|X, W) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

- Training data $D = \{\langle X^1, Y^1 \rangle, \dots, \langle X^L, Y^L \rangle\}$
- Data likelihood = $\prod_l P(X^l, Y^l | W)$
- Data conditional likelihood = $\prod_l P(Y^l | X^l, W)$

$$W_{MCLE} = \arg \max_W \prod_l P(Y^l | W, X^l)$$

Expressing Conditional Log Likelihood

$$l(W) \equiv \ln \prod_l P(Y^l | X^l, W) = \sum_l \ln P(Y^l | X^l, W)$$

$$P(Y = 0 | X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1 | X, W) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$\begin{aligned} l(W) &= \sum_l Y^l \ln P(Y^l = 1 | X^l, W) + (1 - Y^l) \ln P(Y^l = 0 | X^l, W) \\ &= \sum_l Y^l \ln \frac{P(Y^l = 1 | X^l, W)}{P(Y^l = 0 | X^l, W)} + \ln P(Y^l = 0 | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i w_i X_i^l)) \end{aligned}$$

Maximizing Conditional Log Likelihood

$$P(Y = 0|X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1|X, W) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$\begin{aligned} l(W) &\equiv \ln \prod_l P(Y^l | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i^n w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i^n w_i X_i^l)) \end{aligned}$$

Good news: $l(W)$ is concave function of W

Bad news: no closed-form solution to maximize $l(W)$

An example of closed-form solution

The solution of any second-degree polynomial equation can be expressed in terms of its coefficients, using only addition, subtraction, multiplication, division, and **square roots**, in the familiar **quadratic formula**: the roots of the polynomial $ax^2 + bx + c$ (with $a \neq 0$) are

$$\frac{-b \pm \sqrt{b^2 - 4ac}}{2a}.$$

Maximize Conditional Log Likelihood: Gradient Ascent

$$\begin{aligned} l(W) &\equiv \ln \prod_l P(Y^l | X^l, W) \\ &= \sum_l Y^l \left(w_0 + \sum_i w_i X_i^l \right) - \ln \left(1 + \exp \left(w_0 + \sum_i w_i X_i^l \right) \right) \\ \frac{\partial l(W)}{\partial w_i} &= \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W)) \end{aligned}$$


Gradient ascent algorithm: iterate until change $< \varepsilon$
For all i , repeat

$$w_i \leftarrow w_i + \eta \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

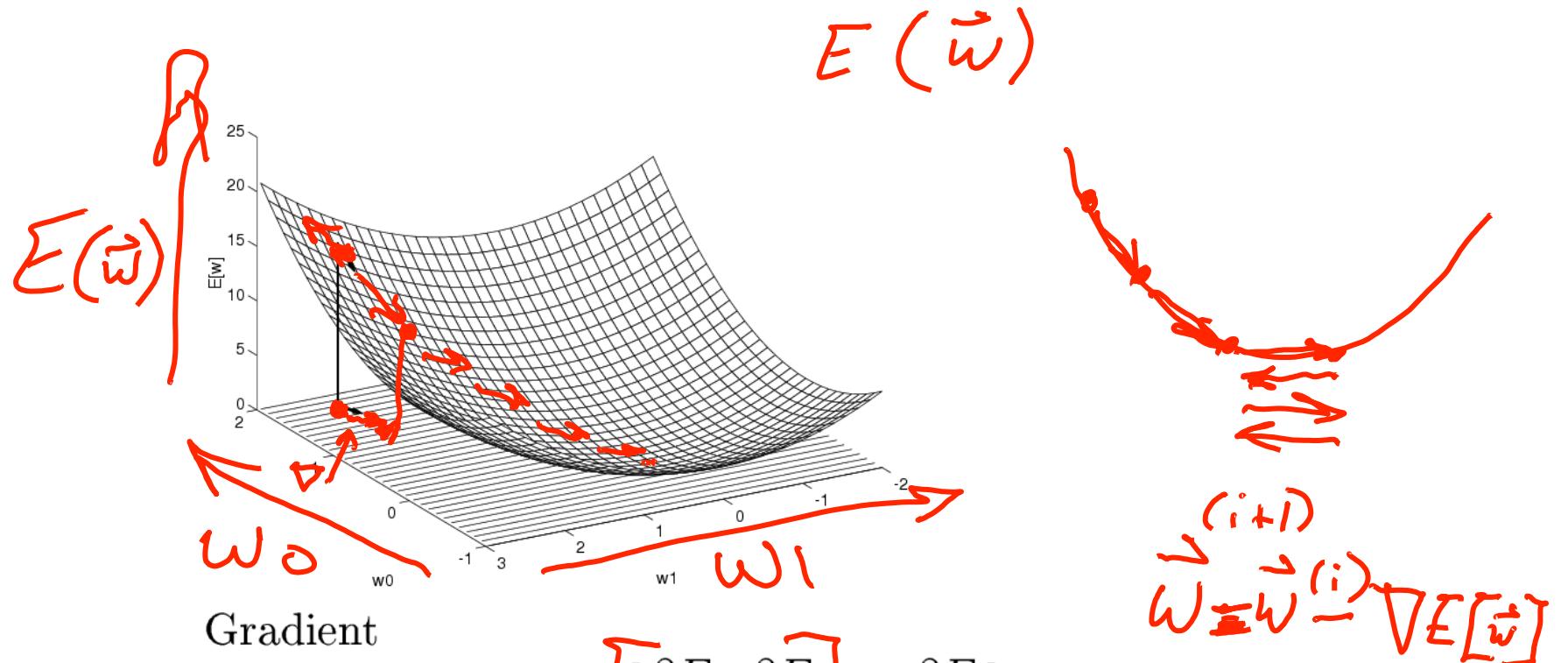
If $f(x) = e^x$, then $f'(x) = e^x$.

If $f(x) = a^x$, $a > 0$, $a \neq 1$, then $f'(x) = (\ln a) \cdot a^x$.

If $f(x) = \ln x$, then $f'(x) = \frac{1}{x}$.

If $f(x) = \log_a x$, $a > 0$, $a \neq 1$, then $f'(x) = \frac{1}{(\ln a) \cdot x}$.

Gradient Descent



$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

i.e.,

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

MAP estimates and Regularization

- Maximum a posteriori estimate with prior $W \sim N(0, \sigma I)$

$$W \leftarrow \arg \max_W \ln [P(W) \prod_l P(Y^l | X^l, W)]$$

$$w_i \leftarrow w_i - \eta \lambda w_i + \eta \sum_l X_i^l (Y^l - \hat{P}(Y^l = 1 | X^l, W))$$

called a “regularization” term

- helps reduce overfitting
- keep weights nearer to zero (if $P(W)$ is zero mean Gaussian prior), or whatever the prior suggests
- used very frequently in Logistic Regression

G.Naïve Bayes vs. Logistic Regression

[Ng & Jordan, 2002]

Recall two assumptions deriving form of LR from GNB:

1. X_i conditionally independent of X_k given Y
2. $P(X_i | Y = y_k) = N(\mu_{ik}, \sigma_i)$, \leftarrow not $N(\mu_{ik}, \sigma_{ik})$

Consider three learning methods:

- GNB (assumption 1 only) -- decision surface can be non-linear
- GNB2 (assumption 1 and 2) – decision surface linear
- LR -- decision surface linear, trained without assumption 1.

Which method works better if we have *infinite* training data, and...

- Both (1) and (2) are satisfied: $LR = GNB2 = GNB$
- (1) is satisfied, but not (2) : GNB > GNB2, GNB > LR, LR > GNB2
- Neither (1) nor (2) is satisfied: GNB > GNB2, LR > GNB2, LR > <GNB

**Use log likelihood to avoid underflow:
prod → sum**

**how can we know if the step size is
reasonable?**

Bias, Variance and Error

Bias and Variance

given algorithm that outputs estimate $\hat{\theta}$ for θ , we define:

the bias of the estimator: $E[\hat{\theta}] - \theta$

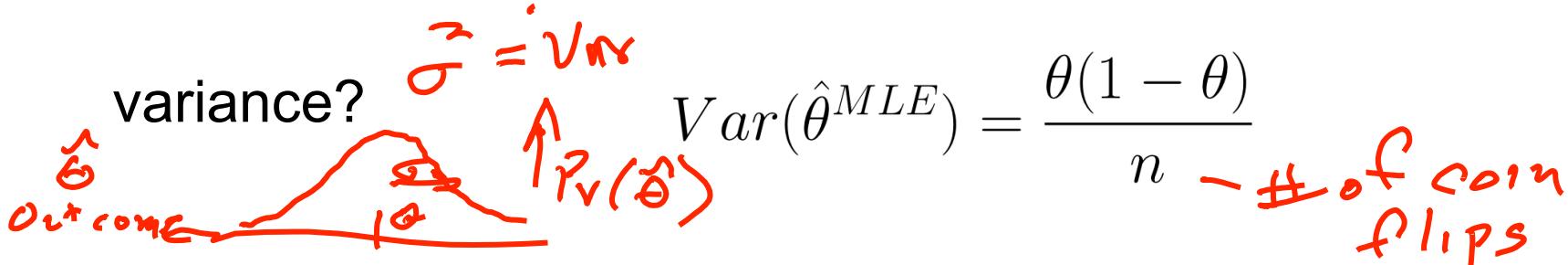
the variance of estimator: $E[(\hat{\theta} - E[\hat{\theta}])^2]$

e.g. $\hat{\theta}^{MLE}$ estimator for probability θ of heads, based on n independent coin flips

$$\hat{\theta}^{MLE} = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

what is its bias?

$$\textcircled{1} \text{ for } \hat{\theta}^{MLE}$$



Here, the p is theta

$$L(p) = \prod_{i=1}^n p^{x_i} (1-p)^{(1-x_i)}$$

$$\ell(p) = \log p \sum_{i=1}^n x_i + \log(1-p) \sum_{i=1}^n (1-x_i)$$

$$\frac{\partial \ell(p)}{\partial p} = \frac{\sum_{i=1}^n x_i}{p} - \frac{\sum_{i=1}^n (1-x_i)}{1-p} \stackrel{\text{set}}{=} 0$$

$$\sum_{i=1}^n x_i - p \sum_{i=1}^n x_i = p \sum_{i=1}^n (1-x_i)$$

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

Bias and Variance $E(\hat{\theta} - \theta)^2$

given algorithm that outputs estimate $\hat{\theta}$ for θ , we define:

the bias of the estimator: $E[\hat{\theta}] - \theta$

the variance of estimator: $E[(\hat{\theta} - E[\hat{\theta}])^2]$

which estimator has higher bias? higher variance?

$$\hat{\theta}^{MLE} = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

$\hat{\theta}_n^{MLE}$ bias = 0 $\text{Var} = \frac{(1-\theta)\theta}{n}$

$$\hat{\theta}^{MAP} = \frac{\alpha_1 + \beta_1 - 1}{(\alpha_1 + \beta_1 - 1) + (\alpha_0 + \beta_0 - 1)}$$

$\hat{\theta}_n^{MAP}$ bias > 0 for finite n
var is less

$$P(|X - E[X]| \geq t) = \frac{\text{Var}(X)}{t^2}$$

$$\mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}(\hat{\theta}^2 + \theta^2 - 2\theta \cdot \hat{\theta}) = \mathbb{E}(\hat{\theta}^2) + \theta^2 - 2\theta \cdot \mathbb{E}(\hat{\theta})$$

remember we have: $\mathbb{E}(\hat{\theta}^2) = \text{Var}(\hat{\theta}) + [\mathbb{E}(\hat{\theta})]^2$

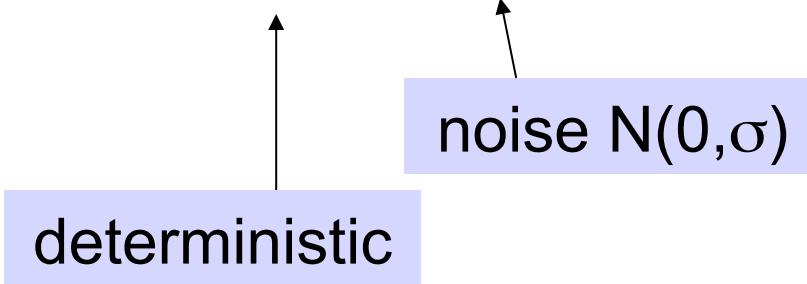
$$\text{So therefore: } \mathbb{E}[(\hat{\theta} - \theta)^2] = \text{Var}(\hat{\theta}) + \underbrace{[\mathbb{E}(\hat{\theta})]^2 + \theta^2 - 2\theta \cdot \mathbb{E}(\hat{\theta})}_{\text{this is square of bias}} = \text{Var}(\hat{\theta}) + [\mathbb{E}(\hat{\theta}) - \theta]^2$$

Bias – Variance decomposition of error

Reading: Bishop chapter 9.1, 9.2

- Consider simple regression problem $f:X \rightarrow Y$

$$y = f(x) + \varepsilon \quad h(x) = w_0 + w_1 x,$$



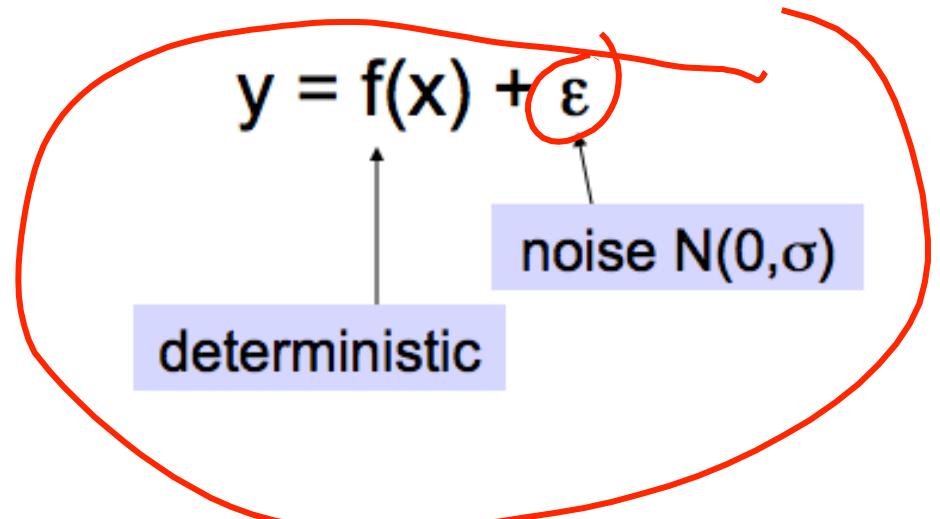
Define the expected prediction error:

$$E_D \left[\int_y \int_x (h(x) - f(x))^2 p(y|x)p(x) dy dx \right]$$

↑
expectation
over
training D

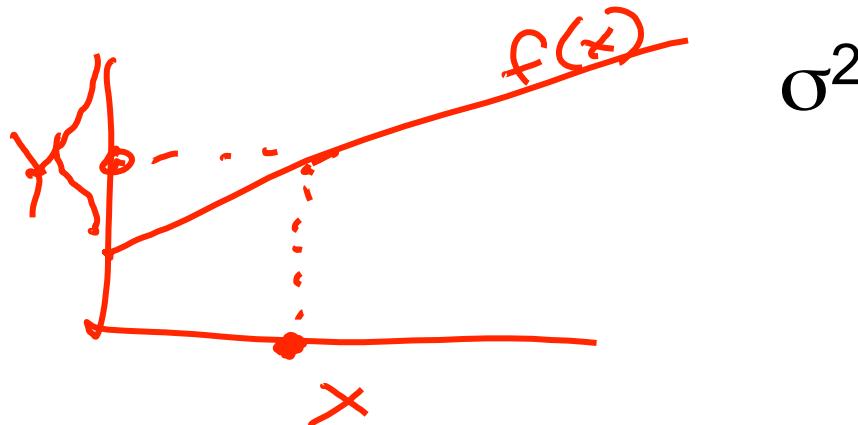
learned
estimate of $f(x)$

Sources of error



What if we have perfect learner, infinite data?

- Our learned $h(x)$ satisfies $h(x)=f(x)$
- Still have remaining, unavoidable error



Sources of error

- What if we have only n training examples?
- What is our expected error
 - Taken over random training sets of size n , drawn from distribution $D=p(x,y)$

$$E_D \left[\int_y \int_x (h(x) - f(x))^2 p(y|x)p(x) dy dx \right]$$

Sources of error

$$y = f(x) + \varepsilon$$

noise $N(0, \sigma)$

deterministic

$$E_D \left[\int_y \int_x (h(x) - \underbrace{f(x)}_2)^2 p(y|x)p(x) dy dx \right]$$
$$= \text{unavoidableError} + \text{bias}^2 + \text{variance}$$

$$\text{bias}^2 = \int (E_D[h(x)] - f(x))^2 p(x) dx$$

$$\text{variance} = \int E_D[(h(x) - E_D[h(x)])^2] p(x) dx$$

Sources of error

$$y = f(x) + \varepsilon$$

noise $N(0, \sigma)$

deterministic

$$E_D \left[\int_y \int_x (h(x) - \underbrace{f(x)}_2)^2 p(y|x)p(x) dy dx \right]$$
$$= \text{unavoidableError} + \text{bias}^2 + \text{variance}$$

$$\text{bias}^2 = \int (E_D[h(x)] - f(x))^2 p(x) dx$$

$$\text{variance} = \int E_D[(h(x) - E_D[h(x)])^2] p(x) dx$$

Graphical Models

- Key Idea:
 - Conditional independence assumptions useful
 - but Naïve Bayes is extreme!
 - Graphical models express sets of conditional independence assumptions via graph structure
 - Graph structure plus associated parameters define *joint probability distribution over set of variables*

Represent Joint Probability Distribution over Variables

Visit to Asia

x_1

Smoking

x_2

Tuberculosis

x_3

Lung Cancer

x_4

Bronchitis

x_5

Tuberculosis
or Cancer

x_6

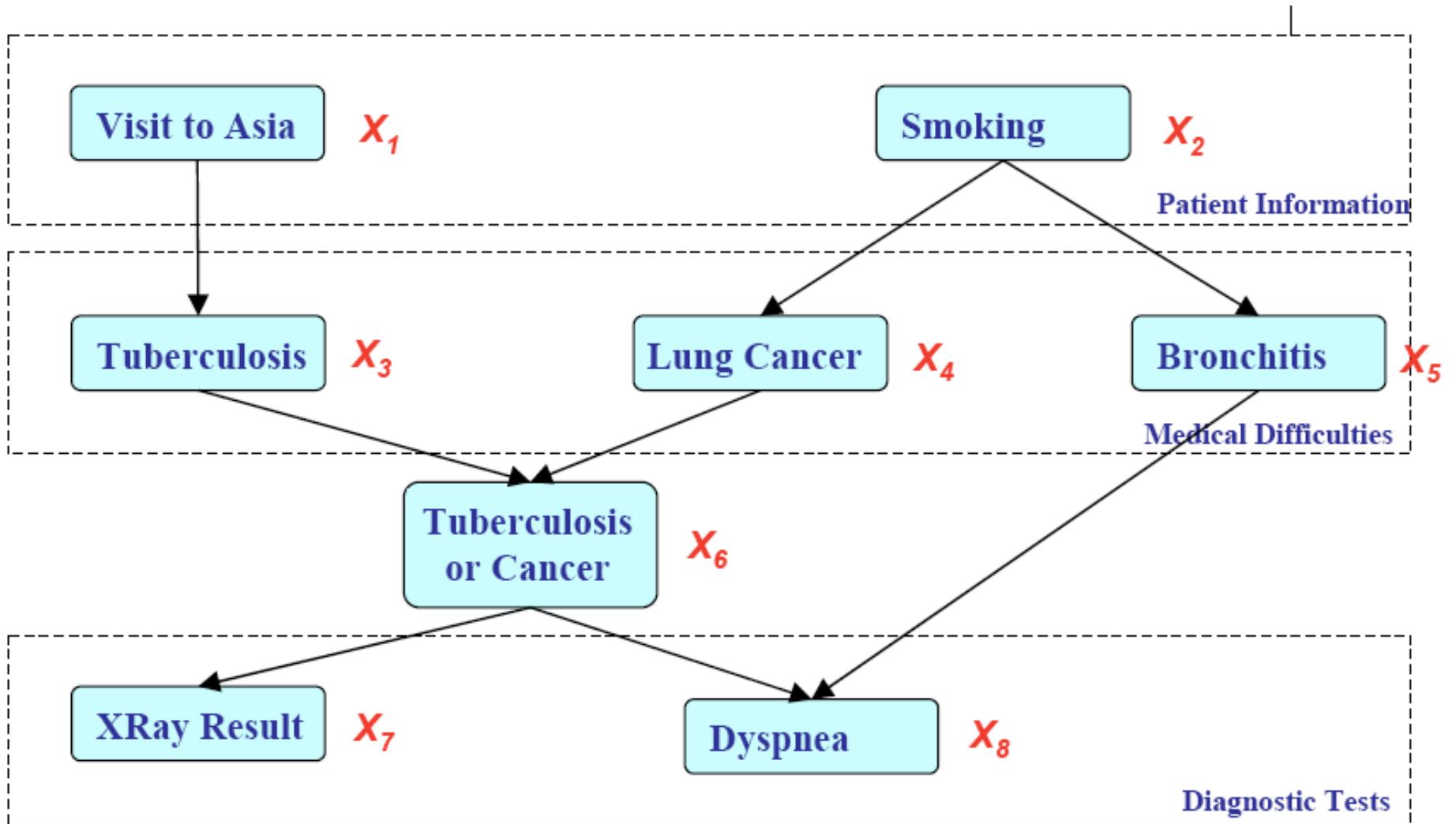
XRay Result

x_7

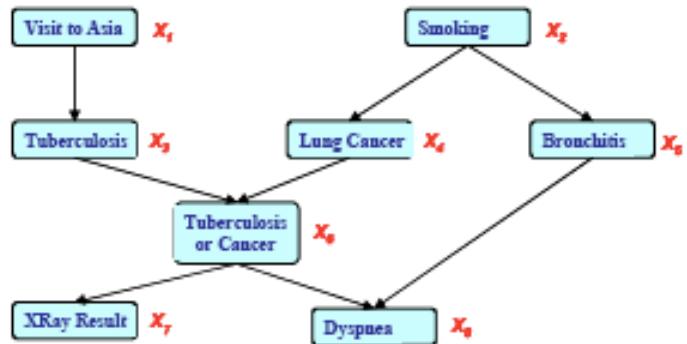
Dyspnea

x_8

Describe network of dependencies



Bayes Nets define Joint Probability Distribution in terms of this graph, plus parameters

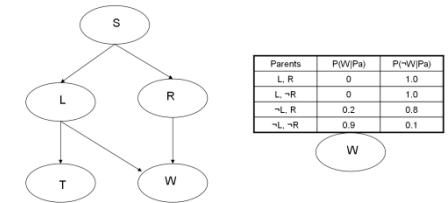


$$\begin{aligned} & P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\ = & P(X_1) P(X_2) P(X_3|X_1) P(X_4|X_2) P(X_5|X_2) \\ & P(X_6|X_3, X_4) P(X_7|X_6) P(X_8|X_5, X_6) \end{aligned}$$

Benefits of Bayes Nets:

- Represent the full joint distribution in fewer parameters, using prior knowledge about dependencies
- Algorithms for inference and learning

Bayesian Networks Definition



A Bayes network represents the joint probability distribution over a collection of random variables

A Bayes network is a directed acyclic graph and a set of conditional probability distributions (CPD's)

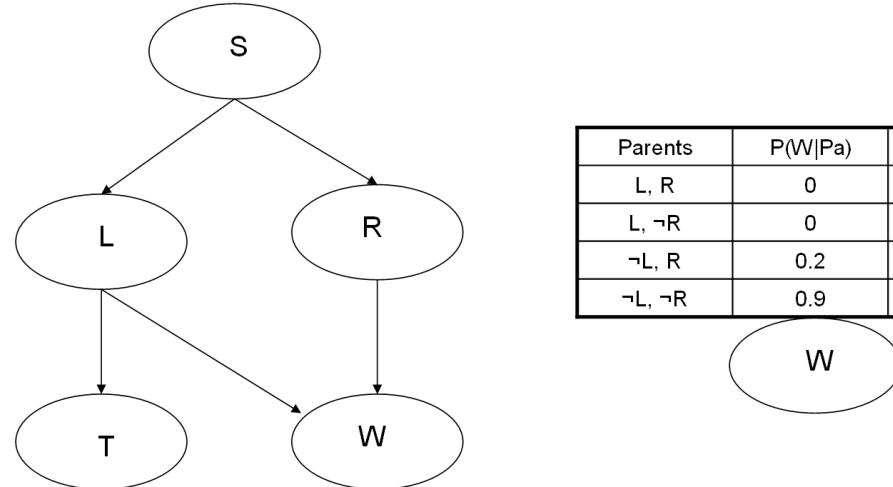
- Each node denotes a random variable
- Edges denote dependencies
- For each node X_i its CPD defines $P(X_i | Pa(X_i))$
- The joint distribution over all variables is defined to be

$$P(X_1 \dots X_n) = \prod_i P(X_i | Pa(X_i))$$

$Pa(X)$ = immediate parents of X in the graph

Bayesian Networks

- CPD for each node X_i describes $P(X_i | Pa(X_i))$



Chain rule of probability says that in general:

$$P(S, L, R, T, W) = P(S)P(L|S)P(R|S, L)P(T|S, L, R)P(W|S, L, R, T)$$

But in a Bayes net: $P(X_1 \dots X_n) = \prod_i P(X_i | Pa(X_i))$

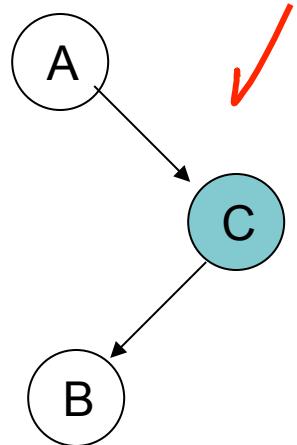
D-Separation: A formal way to analyze conditional independency

Easy Network 1: Head to Tail

prove A cond indep of B given C?

i.e., $p(a,b|c) = p(a|c) p(b|c)$

$$P(A=a) \\ \text{simply } P(a)$$



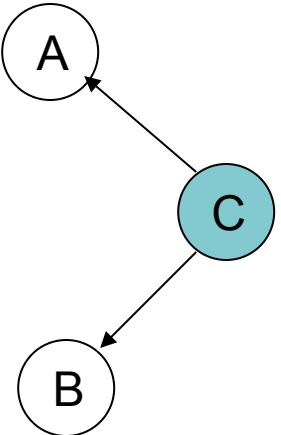
$$p(a,b|c) = p(a|c) p(b|c) \quad \leftarrow A \perp B | c$$

$$p(a,b|c) = \frac{p(a,b,c)}{p(c)} = \frac{p(a) p(c|a) p(b|c)}{p(c)}$$
$$\frac{p(a,c)}{p(c)} = p(a|c)$$

let's use $p(a,b)$ as shorthand for $p(A=a, B=b)$

Easy Network 2: Tail to Tail

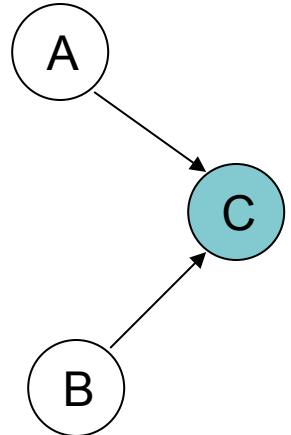
prove A cond indep of B given C? ie., $p(a,b|c) = p(a|c) p(b|c)$



let's use $p(a,b)$ as shorthand for $p(A=a, B=b)$

Easy Network 3: Head to Head

prove A cond indep of B given C? NO!



Summary:

- $p(a,b) = p(a)p(b)$
- $p(a,b|c) \neq p(a|c)p(b|c)$

Explaining away.

e.g.,

- A=earthquake
- B=breakIn
- C=motionAlarm

X and Y are conditionally independent given Z,
if and only if X and Y are D-separated by Z.

[Bishop, 8.2.2]

Suppose we have three sets of random variables: X, Y and Z

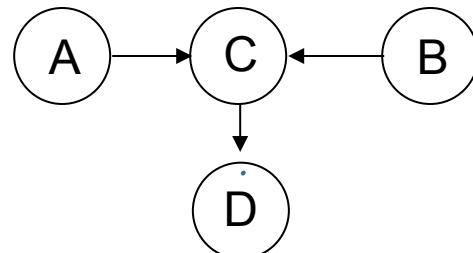
X and Y are D-separated by Z (and therefore conditionally indep, given Z) iff every path from every variable in X to every variable in Y is blocked

A path from variable X to variable Y is **blocked** if it includes a node in Z such that either



1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z

2. or, the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z



X and Y are D-separated by Z (and therefore conditionally indep, given Z) iff every path from any variable in X to any variable in Y is blocked by Z

A path from variable A to variable B is **blocked** by Z if it includes a node such that either

1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z

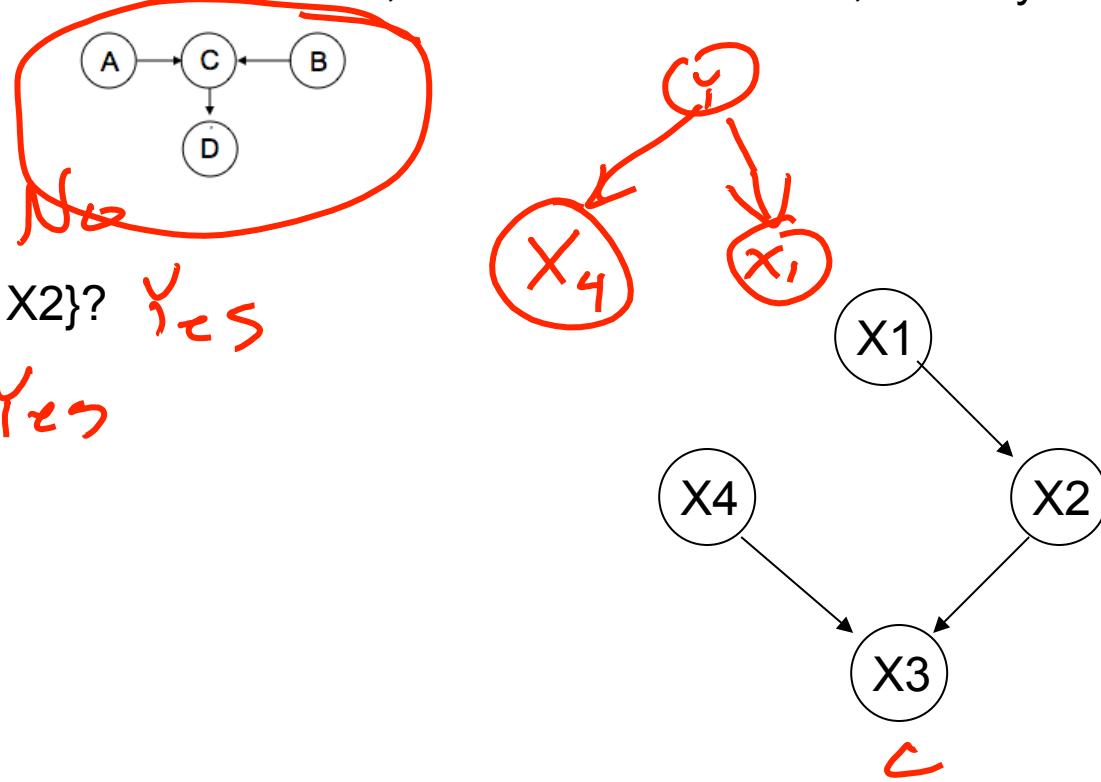


2. the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z

X_4 indep of X_1 given X_3 ? *No*

X_4 indep of X_1 given $\{X_3, X_2\}$? *Yes*

X_4 indep of X_1 given $\{\}$? *Yes*

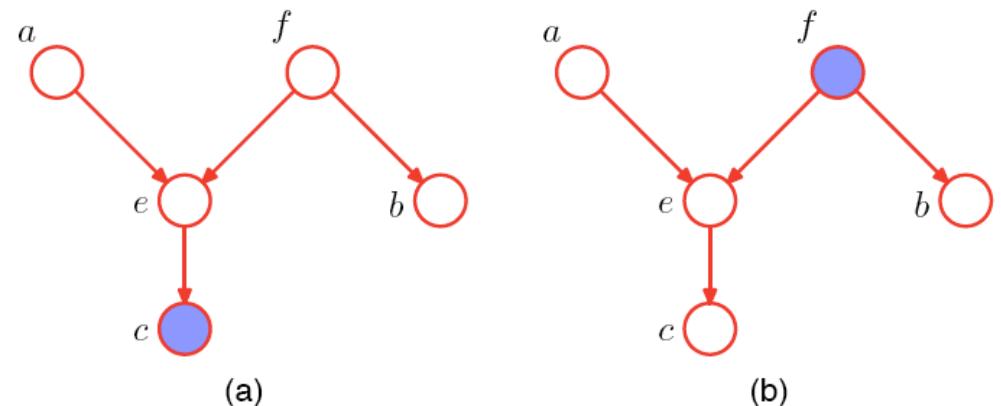


X and Y are **D-separated** by Z (and therefore conditionally indep, given Z) iff every path from any variable in X to any variable in Y is **blocked**

A path from variable A to variable B is **blocked** if it includes a node such that either

1. arrows on the path meet either head-to-tail or tail-to-tail at the node and this node is in Z
2. or, the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in Z

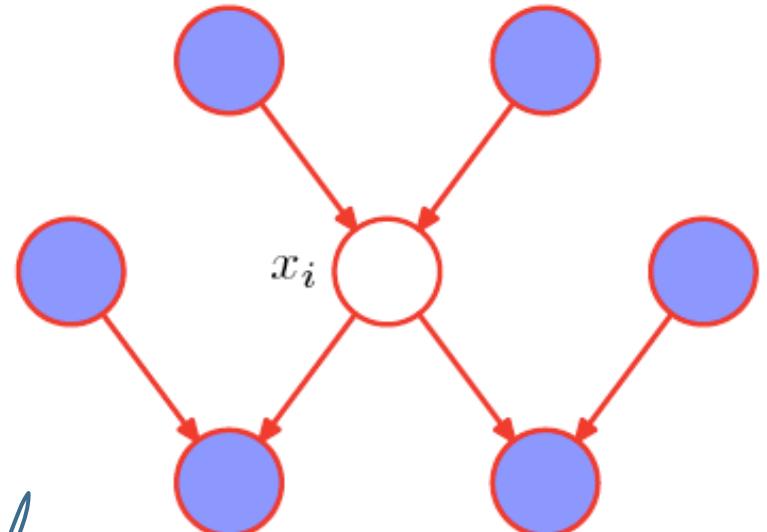
a indep of b given c?



a indep of b given f ?

Markov Blanket

The Markov blanket of a node x_i comprises the set of parents, children and co-parents of the node. It has the property that the conditional distribution of x_i , conditioned on all the remaining variables in the graph, is dependent only on the variables in the Markov blanket.



CO-parent = other side
of x_i 's colliders

from [Bishop, 8.2]

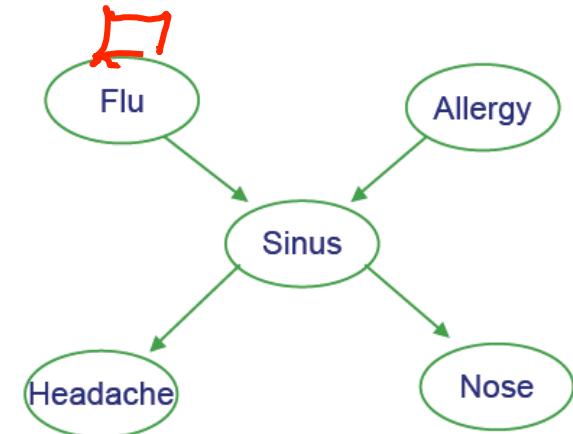
Inference in Bayes Nets

- In general, intractable (NP-complete)
- For certain cases, tractable
 - Assigning probability to fully observed set of variables
 - Or if just one variable unobserved
 - Or for singly connected graphs (ie., no undirected loops)
 - Variable elimination
 - Belief propagation
- Often use Monte Carlo methods
 - e.g., Generate many samples according to the Bayes Net distribution, then count up the results
 - Gibbs sampling
- Variational methods for tractable approximate solutions

see Graphical Models course 10-708

Prob. of joint assignment: easy

- Suppose we are interested in joint assignment $\langle F=f, A=a, S=s, H=h, N=n \rangle$



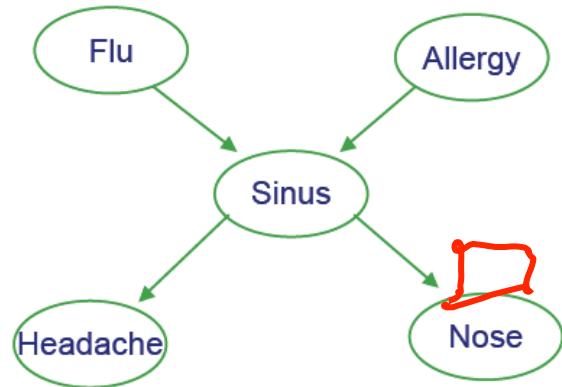
What is $P(f,a,s,h,n)$?

$$P(f) P(a) P(s|f, a) P(h|s) P(n|s)$$

let's use $p(a,b)$ as shorthand for $p(A=a, B=b)$

Prob. of marginals: not so easy

- How do we calculate $P(N=n)$?



$$P(N=n) = \sum_{f, a, h, s} P(N=n, f, a, h, s)$$

$f \in F, a \in A, \dots$

we can:

- 1) sum up the other variables
- 2) do MCMC sampling

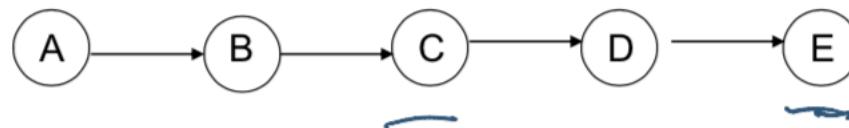
let's use $p(a,b)$ as shorthand for $p(A=a, B=b)$

Let's do $P(C=1)$

Prob. of marginals: not so easy

But sometimes the structure of the network allows us to be clever → avoid exponential work

e.g., chain



Mouse Select Text Draw Stamp Spotlight Eraser Format Undo Redo Clear Save

Let's do $P(C=1)$

Prob. of marginals: not so easy

But sometimes the structure of the network allows us to be clever → avoid exponential work



(a, b, c, d, e)

~~H(C = 1)
ex. chain
P(C = 1) =~~

A B C D E $\rightarrow 2^4$

$$P(C=1) = \sum P(a b c=1 | d e)$$

~~abde~~

$$2^2 = \sum_{abde} P(a) P(b|a) P(c|b) P(d|c) P(e|d)$$

$$= \sum_{abde} P(ab) P(c|ab) P(de|c=1)$$

$$= \sum_{ab} P(ab) P(c=1|ab) \cancel{\sum P(de|c=1)}$$

From Yunjie(Chloe) Song to Everyone:
that node is a co-parent?

From Patrick to Everyone:
yes

From Yunjie(Chloe) Song to Everyone:
not in frequentist

From Patrick to Everyone:
1?

From Ruhao Xin to Everyone:
1?

From Patrick to Everyone:
but since a and c are conditional given B, do we still need to cons

To: Yuqin Bai (Privately)

Type message here...

Learning of Bayes Nets

- Four categories of learning problems
 - Graph structure may be known/unknown
 - Variable values may be fully observed / partly unobserved
- Easy case: learn parameters for graph structure is *known*, and data is *fully observed*
- Interesting case: graph *known*, data *partly known*
- Gruesome case: graph structure *unknown*, data *partly unobserved*

Learning CPTs from Fully Observed Data

- Example: Consider learning the parameter

$$\theta_{s|ij} \equiv P(S = 1 | F = i, A = j)$$

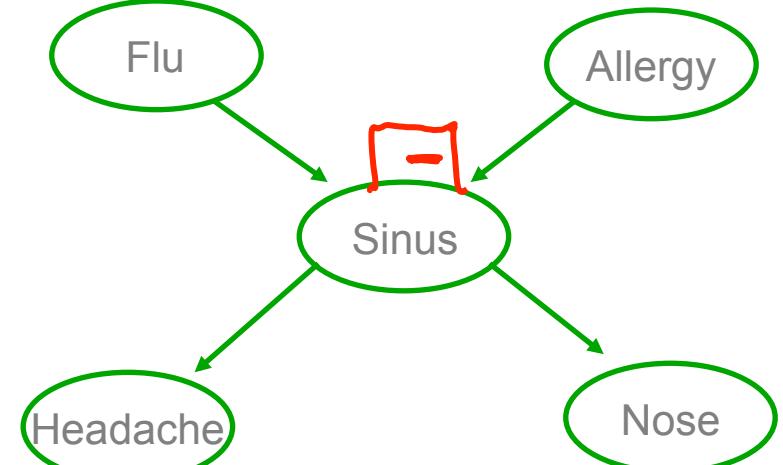
- Max Likelihood Estimate is

$$\theta_{s|ij} = \frac{\sum_{k=1}^K \delta(f_k = i, a_k = j, s_k = 1)}{\sum_{k=1}^K \delta(f_k = i, a_k = j)}$$

kth training example

$\delta(x) = 1$ if $x=\text{true}$,
 $= 0$ if $x=\text{false}$

- Remember why?



let's use $p(a,b)$ as shorthand for $p(A=a, B=b)$

MLE estimate of $\theta_{s|ij}$ from fully observed data

- Maximum likelihood estimate

$$\theta \leftarrow \arg \max_{\theta} \log P(\text{data}|\theta)$$

- Our case:

$$P(\text{data}|\theta) = \prod_{k=1}^K P(f_k, a_k, s_k, h_k, n_k)$$

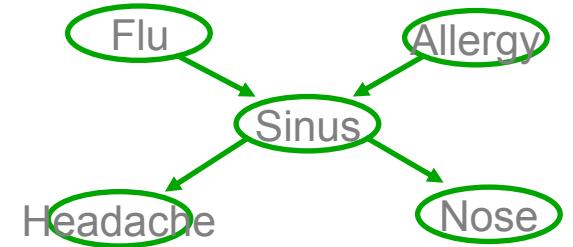
↙+true example ↘

$$P(\text{data}|\theta) = \prod_{k=1}^K P(f_k)P(a_k)P(s_k|f_k a_k)P(h_k|s_k)P(n_k|s_k)$$

$$\log P(\text{data}|\theta) = \sum_{k=1}^K \log P(f_k) + \log P(a_k) + \log P(s_k|f_k a_k) + \log P(h_k|s_k) + \log P(n_k|s_k)$$

$$\frac{\partial \log P(\text{data}|\theta)}{\partial \theta_{s|ij}} = \sum_{k=1}^K \frac{\partial \log P(s_k|f_k a_k)}{\partial \theta_{s|ij}}$$

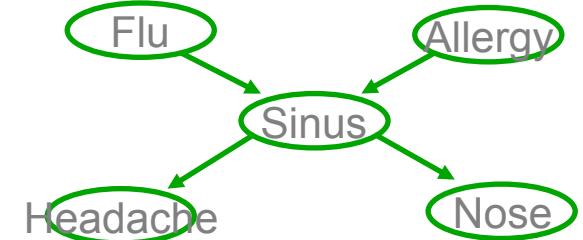
$$\theta_{s|ij} = \frac{\sum_{k=1}^K \delta(f_k = i, a_k = j, s_k = 1)}{\sum_{k=1}^K \delta(f_k = i, a_k = j)}$$



Estimate θ from partly observed data

- What if FAHN observed, but not S?
- Can't calculate MLE

$$\theta \leftarrow \arg \max_{\theta} \log \prod_k P(f_k, a_k, s_k, h_k, n_k | \theta)$$



- Let X be all *observed* variable values (over all examples)

- Let Z be all *unobserved* variable values

- Can't calculate MLE:

$$\theta \leftarrow \arg \max_{\theta} \log P(X, Z | \theta)$$

- WHAT TO DO?

EM Algorithm - Informally

EM is a general procedure for learning from partly observed data

Given observed variables X, unobserved Z ($X=\{F,A,H,N\}$, $Z=\{S\}$)

Begin with arbitrary choice for parameters θ

Iterate until convergence:

- E Step: estimate the values of unobserved Z, using θ
- M Step: use observed values plus E-step estimates to derive a better θ

Guaranteed to find local maximum.

Each iteration increases $E_{P(Z|X,\theta)}[\log P(X, Z|\theta')]$

EM Algorithm - Precisely

EM is a general procedure for learning from partly observed data

Given observed variables X, unobserved Z ($X=\{F,A,H,N\}$, $Z=\{S\}$) ✓

Define $Q(\theta'|\theta) = E_{P(Z|X,\theta)}[\log P(X, Z|\theta')]$

w^{et} *current* *current* *M step new*

Iterate until convergence:

- E Step: Use X and current θ to calculate $P(Z|X,\theta)$
- M Step: Replace current θ by

$$\theta \leftarrow \arg \max_{\theta'} Q(\theta'|\theta)$$

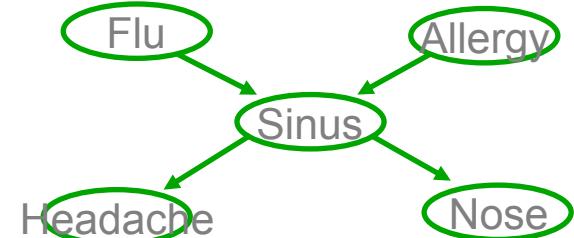
Guaranteed to find local maximum.

Each iteration increases $E_{P(Z|X,\theta)}[\log P(X, Z|\theta')]$

back to this example, say S is Z .

- EM seeks estimate:

$$\theta \leftarrow \arg \max_{\theta} E_{Z|X,\theta} [\log P(X, Z|\theta)]$$



- here, observed $X=\{F,A,H,N\}$, unobserved $Z=\{S\}$

$$\log P(X, Z|\theta) = \sum_{k=1}^K \log P(f_k) + \log P(a_k) + \log P(s_k|f_k a_k) + \log P(h_k|s_k) + \log P(n_k|s_k)$$

$$E_{P(Z|X,\theta)} \log P(X, Z|\theta) = \sum_{k=1}^K \sum_{i=0}^1 P(s_k = i | f_k, a_k, h_k, n_k) = \frac{P(s_k = 0, f_k, a_k, h_k, n_k)}{P(s_k = 0, f_k, a_k, h_k, n_k) + P(s_k = 1, f_k, a_k, h_k, n_k)}$$

\downarrow \downarrow

$$[\log P(f_k) + \log P(a_k) + \log P(s_k | f_k a_k) + \log P(h_k | s_k) + \log P(n_k | s_k)]$$

$$E[x] = \sum_i P(x=i) i$$

coin: $\rightarrow \theta \rightarrow \begin{cases} 0 \\ 1 \end{cases} X_1$
 $Y \rightarrow \begin{cases} 0 \\ 1 \end{cases} X_2$

$$P(Y=1) = \frac{1}{2}$$

\downarrow
 $X = 1 \ 2 \ 3 \ 4 \ 5 \ 6$

(X)

4	5	5	15	15	15
5	5	5	5	5	5
5	5	5	5	5	5

$$\begin{aligned} E(X|Y=1) &= \sum_x x P(X|Y=1) = \\ &= \sum_x x P(X|Y=1) = \end{aligned}$$

