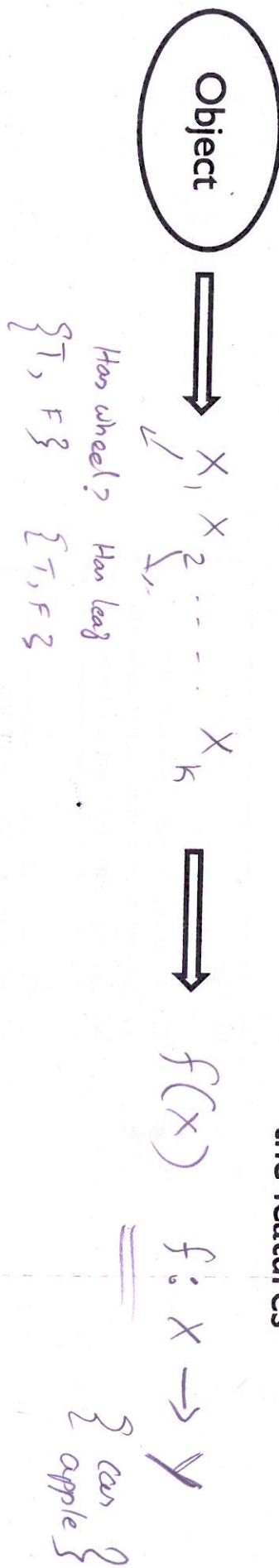


Learning to Predict

Features ✓
Classifier, a function over ✓
the features



Two main questions in ML:

1) How do we get features?

- Domain experts / problem study to give features
- Deep learning allocates learning features too!
from raw input

2) How do we obtain the classifier function?

linear function: $\sum w_i x_i + b$ ↗ linear fun.

non-linear functions

$(x^T x)^{\alpha \beta}$

General Naïve Bayes Model

> Suppose you want to predict if an email is spam or ham.

> What are your outputs?

$$Y = \{\underline{\text{Spam}}, \underline{\text{Ham}}\}$$

What are your features?

1. has-nigerian-prince = $F_1 = \{\tau, F\}$

2. has-\$\\$\\$ = $F_2 = \{\tau, F\}$

3. is-grammatical = F_3

4. ALL-CAPS? = F_4



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret....

TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES FOR ONLY \$99



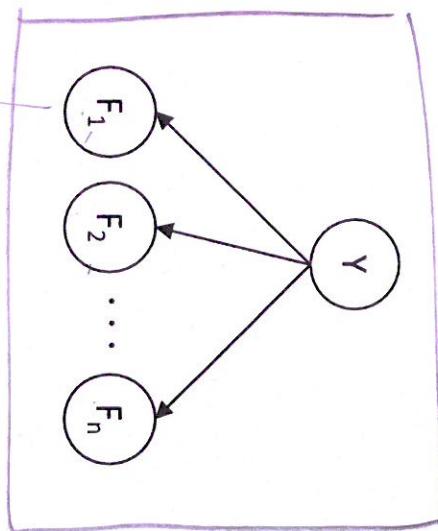
Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use. I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

Let's denote these by random variables F_i .

Naïve Bayes is a Bayesian Network!

- Assume all features are independent effects of label.
 - We only have to specify how each feature depends on the class

$$\begin{aligned}
 P(Y, F_1 \dots F_n) &= p_y(y) \times p_{F_1}(F_1|y) \times p_{F_2}(F_2|y) \\
 &\quad \cdots \times p_{F_n}(F_n|y) \\
 &= p_y(y) \prod_{i=1}^n p_{F_i}(F_i|y)
 \end{aligned}$$



- What are the CPTs associated with this network?

~~Y~~

	$p_y(y)$
--	----------

~~F_1~~

	$p_{F_1}(F_1 y)$
--	------------------

Ham

Spam

Ham

Spam

F

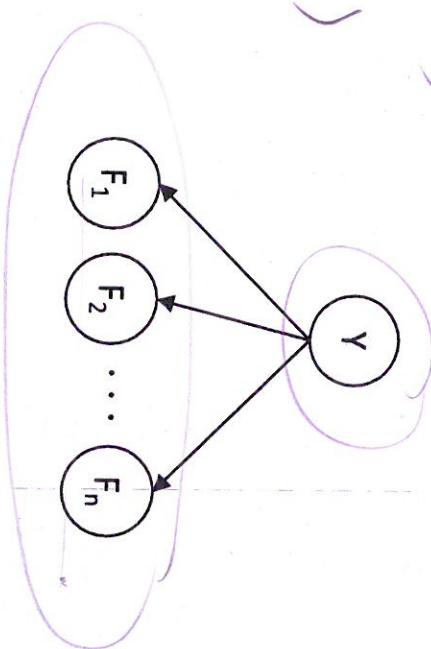
F

Naïve Bayes is a Bayesian Network!

- What probability do we want for classification?

$$P_Y(Y = \text{Spam} \mid F_1 = f_1, F_2 = f_2, \dots, F_n = f_n)$$

$$P_Y(Y = \text{Ham} \mid F_1 = f_1, \dots, F_n = f_n)$$



- How to get the conditional from the joint for this network? Inference!

$$P_Y(Y \mid F_1, \dots, F_n) = \frac{P_Y(Y, F_1, \dots, F_n)}{P_Y(F_1, \dots, F_n)}$$

$$\text{argmax}_{Y \in \{\text{Spam, Ham}\}} P_Y(Y \mid F_1, \dots, F_n) = \text{argmax}_{Y \in \{\text{Spam, Ham}\}} \frac{P_Y(Y, F_1, \dots, F_n)}{P_Y(F_1, \dots, F_n)} = \text{argmax}_Y \prod_{i=1}^n P_Y(Y \mid F_i)$$

Naïve Bayes is a Bayesian Network!

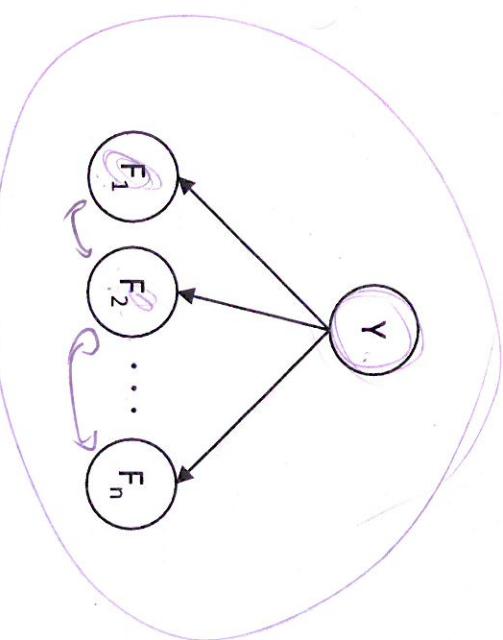
- How do we get the conditional probabilities? Estimate from training data.

$$\cancel{P_{\text{Spam}} = \frac{\#(\text{Spam})}{N}}$$

$$\cancel{P_{\text{Spam} | F_1, F_2, \dots, F_n} = ?}$$

$$P_{\text{Spam}} = \frac{\#(\text{Spam})}{N} \leftarrow \text{total # examples}$$

$$P_{\text{Spam} | F_1, F_2, \dots, F_n} = \frac{\#(F_i = \text{True}, \text{Spam})}{\#(\text{Spam})}$$



An example: Bag-of-words Naïve Bayes

- Data:
 - Collection of emails, labeled spam or ham
 - Note: someone has to hand label all this data!
- Features
 - Words that occur in the document.
- Classifier
 - Bayesian Network
 - Label conditioned on feature variables (spam vs. ham)
 - Assume features are conditionally independent given class label

Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE
MAILINGS, SIMPLY REPLY TO THIS
MESSAGE AND PUT "REMOVE" IN THE
SUBJECT.

99 MILLION EMAIL ADDRESSES
FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

An example: Bag-of-words Naïve Bayes

- BN Variables?

$$y = \{s, H\}$$

w_i = word in position i

e.g.

w_1 = Nigerian

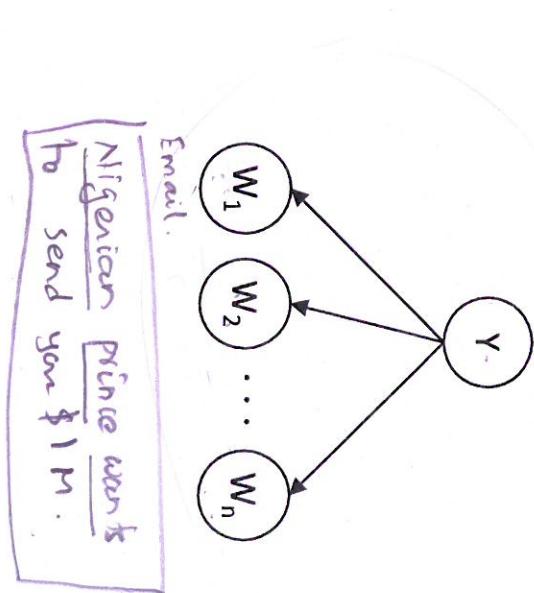
w_2 = Prince

w_3 = wants

$$w_i = \{a, an, the, \dots, zebra\}$$

- BN CPTs?

- ignore for now!



An example: Bag-of-words Naïve Bayes

- Joint Distribution?

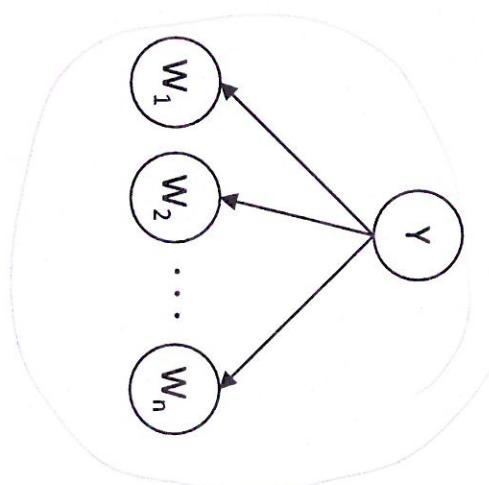
$$P_{\gamma}(y, w_1, \dots, w_n) = P_{\gamma}(y) \prod_{i=1}^n P_{\gamma}(w_i | y)$$

- What do we want for prediction?

$$P_{\gamma}(y | w_1, \dots, w_n)$$

- Use Bayes Rule

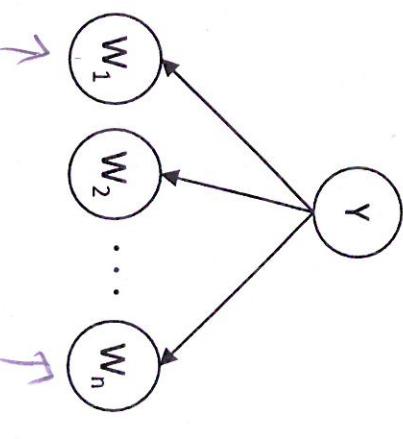
$$\text{argmax}_{\gamma} P_{\gamma}(y) \prod_{i=1}^n P_{\gamma}(w_i | y)$$



An example: Bag-of-words Naïve Bayes

- Use training data to estimate!

$$P(Y) = \frac{\#(Y = \text{spam})}{N} \rightarrow \frac{\#(Y = H)}{N} \rightarrow \# \text{ emails.}$$



$$P(W|\text{spam})$$

$$\Pr(w_i | \text{spam}) \xrightarrow{\# \text{ spam emails in which } w_i \text{ was found}} \frac{\#(w_i = \text{spam})}{\# \text{ spam emails.}}$$

$\Pr(w_1 = \text{Niger} | \text{spam}) = \Pr(w_n = \text{Niger} | \text{spam})$

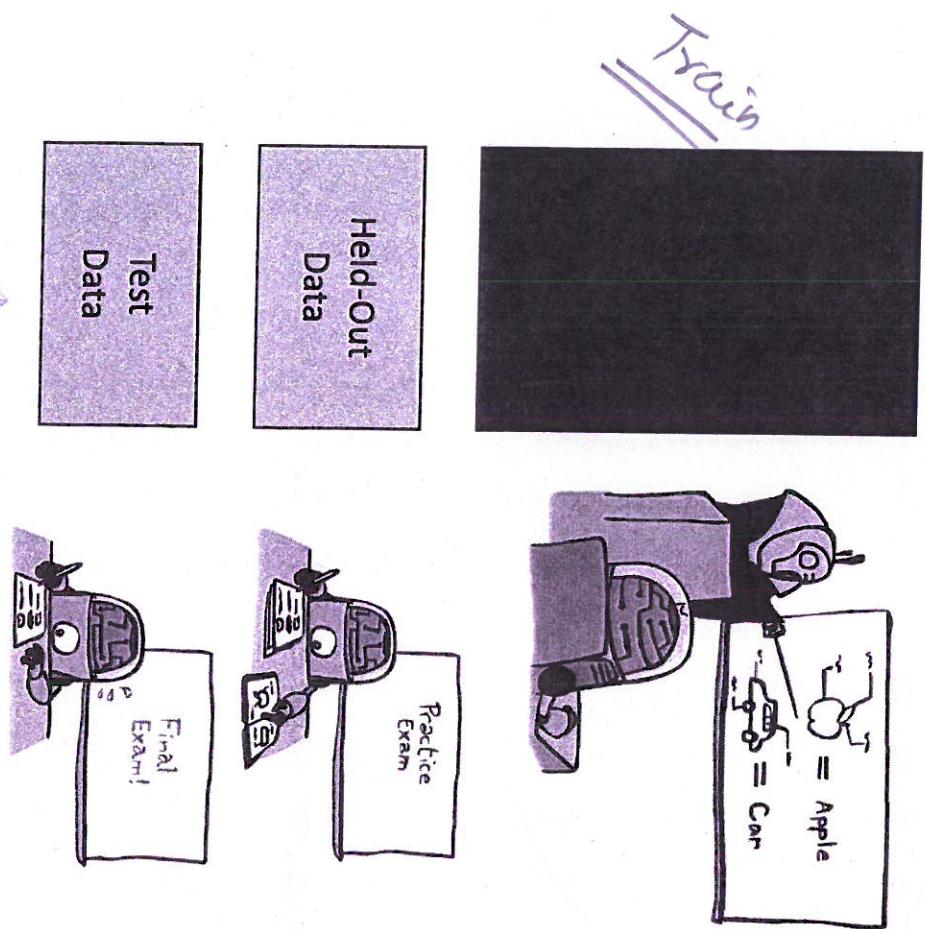
- Maximum likelihood estimation

- Max a posteriori estimation \rightarrow expectations on how the probabilities should look like.

MAP

How do we know our classifier is any good?

Experimental Cycle



Evaluation Metrics:

$$\text{Accuracy} = \frac{\#(\text{correct labels})}{\#(\text{labels})}$$

$$\text{Precision} = \frac{\#(\text{model} = \text{spam}, \text{truth} = \text{spam})}{\#(\text{model} = \text{spam})}$$

multiple
writers

$$\text{Recall} = \frac{\#(\text{model} = s, \text{truth} = s)}{\#(\text{truth} = s)}$$

$$\text{F-measure} = \frac{2 \cdot PR}{P+R}$$

Harmonic mean of P & R

Issues: Independence assumptions

- Need more features—words aren't enough!
 - Have other people just gotten the same email?
 - Is the email in ALL CAPS?
 - Do inline URLs point where they say they point?
- Can add these information sources as new variables in the NB model
 - But will be violating independence assumptions even more.
- Example:
 - “New” and “York” are dependent even when the class is observed.
 - Knowing “York” tells me New is likely to be present in the email.
 - If you $P(\text{“New"}|C)$ and $P(\text{“York"}|C)$ both contribute independent evidence.
- There are classifiers which let you easily add arbitrary features more easily.
 - We'll see one such next.

Discriminative Classification

- Bayesian networks are generative models ✓
 - They model the joint distribution of the instances and the labels.
— Generative in the sense that they can gen. samples.
 - Why model joint $P(Y, X)$?
 - After all we only care about $P(Y|X)$
- We are interested in selecting the best Y given the evidence X i.e., $P(Y|X)$
- Models that directly estimate this are called discriminative classifiers!

e.g. Logistic Regression

Discriminative Classification = Conditional Probability Estimation

- Let's learn a function that can predict $P(Y|X)$ directly!
- Assume that $P(Y|X)$ is a parameterized function of some form.
$$P(Y|X) \leq h_{\mathbf{w}}(x, y) \rightarrow h_{\mathbf{w}}(x) = \begin{cases} 1 & \text{if } x \text{ is spam} \\ 0 & \text{if } x \text{ is not spam} \end{cases}$$

For spam email (x^i, y^i) $h_{\mathbf{w}}(x^i, y^i) = P(y = \text{spam} | x^i) = 1$
- Find the "best" values for these parameters using the training data.
- Logistic regression is one example.

$$h_{\mathbf{w}}(x) \rightarrow \text{has free parameters, } \mathbf{w} = (w_1, w_2, \dots, w_n)$$
$$\begin{matrix} f_1 & f_2 & \dots & f_k \end{matrix}$$

Logistic Regression Recipe: For binary classification

- Assume that there is a function that approximates $P(Y|X)$

$$h_w(x)$$

- We will jump a couple of steps to arrive at this form.

- First, what constraints should this function satisfy?

$$h_w(x) \leq P(y|x)$$

$$\Rightarrow 0 \leq h_w(x) \leq 1$$

Training: Setting an objective function

- Select a set of parameters that maximize the probability of the training data.

$$\text{Train} = \{(y_c^1, x_c^1), (y_c^2, x_c^2), \dots, (y_c^L, x_c^L)\}$$

\downarrow
Email 1 Email 2 ...

- Probability of the training data

$$\prod_{i=1}^L P_R(y_c^i | x_c^i)$$

$\prod_{i=1}^L P_R(y_c^i | x_c^i)$
|| \hookrightarrow Prob. of entire
training data
 ~~$P_R(y_c^i | x_c^i) P_R(x_c^i)$~~

- Find argmax to give the best parameter settings.

ignore

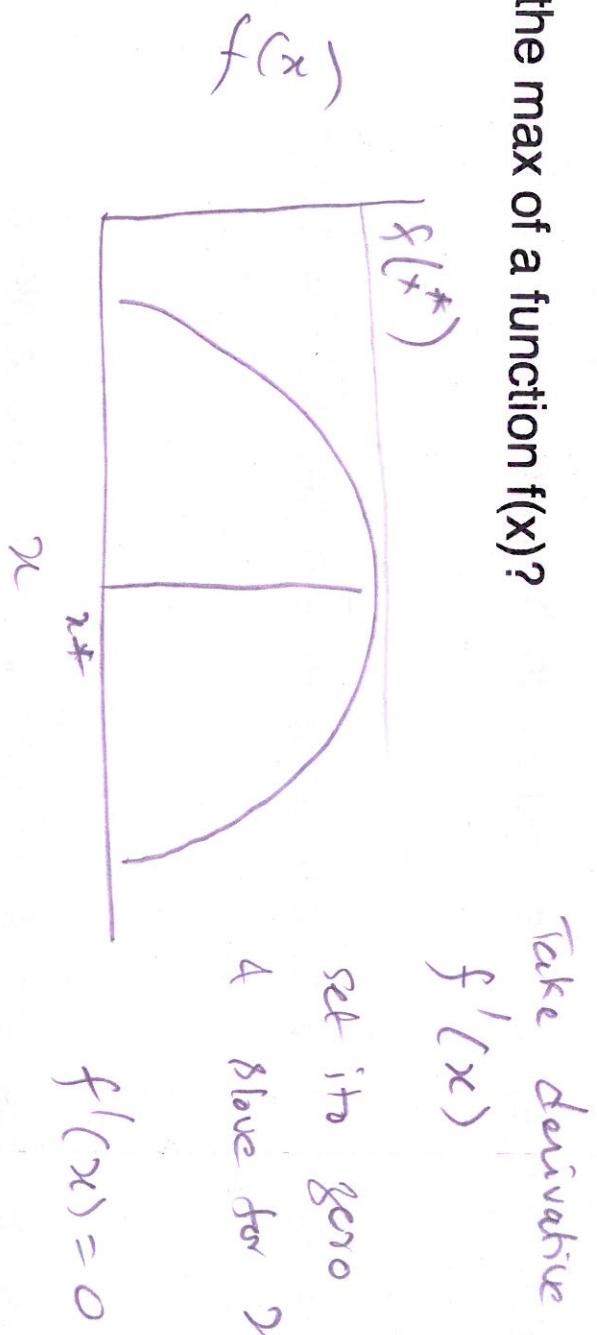
$$\begin{aligned} \text{argmax}_W & \prod_{\ell} P_R(y_c^\ell | x_c^\ell) \\ &= \prod_{\ell} h_W(x_c^\ell) \overline{\prod_{\ell \in \text{ham}} (1 - h_W(x_c^\ell))} \end{aligned}$$

Optimization: Finding the max (or min) of a function

- How to find the *max* of our objective?

$$\text{obj}(\underline{\mathbf{w}}) = \underset{l \in S}{\min} \prod_{i=1}^n h_i(x^l)$$

- How do you find the max of a function $f(x)$?



Optimization: Finding the max (or min) of a function

- Let's compute the derivative of our objective and see if we can solve for W .
$$\text{Obj}(w) = \prod_l P_{\pi}(y^l = y_c^l | x^l) = h_w$$
- Let's make our lives simpler by using a log-likelihood objective instead of likelihood.

- Taking log leaves the argmax unchanged.

$$\underset{w}{\operatorname{argmax}} \sum_l \ln \left\{ P_{\pi}(y^l = \underline{y_c^l} | x^l) \right\}$$

- Rewrite the summation to separate cases w/ $\underline{y^l = 1}$ and $\underline{y^l = 0}$

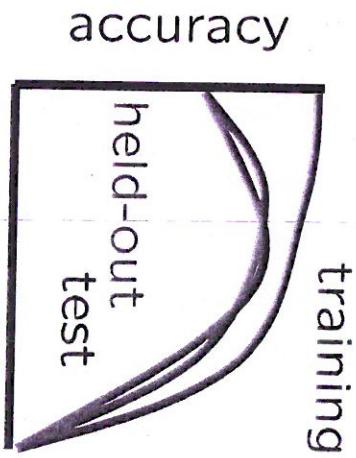
$$\underset{w}{\operatorname{argmax}} \sum_l y_c^l \ln \underline{P_{\pi}(y^l = 1 | x^l)} + (1 - y_c^l) \ln \underline{P_{\pi}(y^l = 0 | x^l)}$$

- Rearrange the terms to collect y_c terms in one place

$$\underset{w}{\operatorname{argmax}} \sum_l y_c^l \ln \left\{ P_{\pi}(y^l = 1 | x^l) \right\} + \ln P_{\pi}(y^l = 0 | x^l)$$

Issues: Generalization and Overfitting

- Relative frequency parameters will overfit the training data!
 - Unlikely that every occurrence of “minute” is 100% spam
 - Unlikely that every occurrence of “seriously” is 100% ham
 - What about all the words that don’t occur in the training set at all?
 - In general, we can’t go around giving unseen events zero probability
- As an extreme case, imagine using the entire email as the only feature
 - Would get the training data perfect (if deterministic labeling)
 - Wouldn’t generalize at all
 - Just making the bag-of-words assumption gives us some generalization, but isn’t enough
- To generalize better: we need to smooth or regularize the estimates
 - Laplace’s estimate
 - Pretend you saw every outcome k extra times



- Two kinds of unknowns, parameters and hyperparameters.
 - Estimate probabilities from training data.
 - Choose k based on held-out performance.

Issues: Independence assumptions

- Need more features—words aren't enough!
 - All caps
 - email domain
- Can add these information sources as new variables in the NB model
 - But the features are not independent -
- Example:
 - "New" and "York"
 - Knowing "York" tells me New occurred before York
 - If you use $P(\text{"New"}|C)$ and $P(\text{"York"}|C)$ we are double counting.
- There are classifiers which let you easily add arbitrary features more easily.
 - We'll see one such next.

Logistic Regression Recipe: For binary classification

- Let's consider a linear combination function.

- Something that simply does a weighted combination of the feature values.

Suppose $h_w(x) = w^T x = \sum_{i=1}^n w_i x_i$

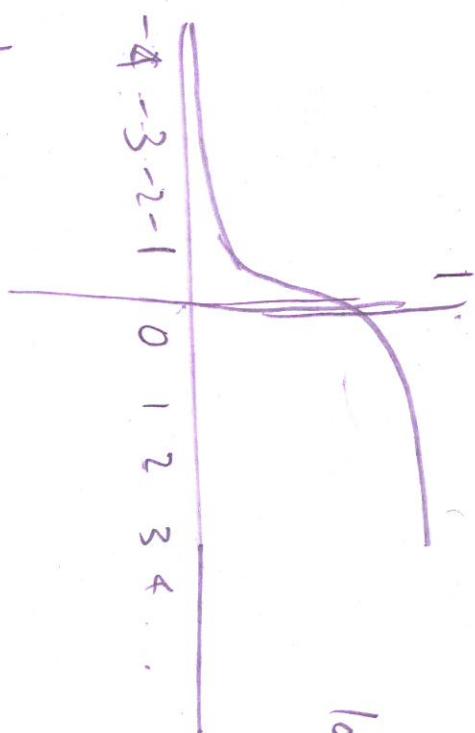
- Does this form satisfy our constraints?

$$0 \leq h_w(x) \leq 1 \quad w_i \in \mathbb{R} \quad \forall w_i \in \mathbb{R}$$

- What is the fix?

$$h_w(x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

logistic function.



If $h_w(x) < 0.5 \Rightarrow \text{ham}$
 $h_w(x) \geq 0.5 \Rightarrow \text{spam}$.

Optimization: Finding the max (or min) of a function

- Let's compute the derivative of our objective and see if we can solve for W .

- Let's make our lives simpler by using a log-likelihood objective instead of likelihood.

- Taking log leaves the argmax unchanged.

$$L(\underline{w})$$

$$\operatorname{argmax}_W \sum_l \ln P(Y^l = y_c^l | X^l) \quad \xrightarrow{\text{log-likelihood objective.}} \quad L(\underline{w})$$

- Rewrite the summation to separate cases w/ $Y^l = 1$ and $Y^l = 0$

$$\operatorname{argmax}_W \sum_l y_c^l \ln \{P(Y^l = 1 | X^l)\} + (1 - y_c^l) \ln \{P(Y^l = 0 | X^l)\}$$

- Rearrange the terms to collect y_c terms in one place

$$\operatorname{argmax}_W \sum_l y_c^l \ln \frac{P(Y^l = 1 | X^l)}{P(Y^l = 0 | X^l)} + \ln P(Y^l = 0 | X^l)$$

Optimization: Find the max of the log-likelihood function.

- Doing a bit of arithmetic will get you
- $$h_w(x) = \frac{e^{w^T x}}{1 + e^{w^T x}} \quad P(y^l = 0 | x^l) = 1 - h_w(x^l) = 1 - \frac{e^{w^T x}}{1 + e^{w^T x}} = \frac{1}{1 + e^{w^T x}}$$

$$\operatorname{argmax}_W \sum_l y_c^l \ln \frac{P(Y^l = 1 | X^l)}{P(Y^l = 0 | X^l)} + \ln P(Y^l = 0 | X^l)$$

$$y_c^l \ln(e^z) - y_c^l z \quad \begin{matrix} \cancel{\text{if } e^z} \\ \cancel{\text{if } 1 + e^z / 1/e^z} \end{matrix}$$

$$W^T X = \sum_{i=1}^n w_i X_i$$

$$\frac{1}{1 + e^{W^T X}}$$

$$w^T x = z$$

- Substituting back

$$\operatorname{argmax}_W \sum y_c^l w^T x + \ln \left(\frac{1}{1 + e^{w^T x}} \right)$$

- Using $\ln(1/a) = -\ln(a) \Rightarrow$

$$\operatorname{argmax}_W \sum y_c^l w^T x - \ln \left(1 + e^{w^T x} \right)$$

Optimization: Find the max of the log-likelihood function.

- Doing a bit of arithmetic will get you

$\xrightarrow{\text{Sub } h_w}$

$$\operatorname{argmax}_W \sum_l y_c^l \ln \frac{P(Y^l = 1 | X^l)}{P(Y^l = 0 | X^l)} + \ln P(Y^l = 0 | X^l)$$



$$W^T X = \sum_{i=1}^n w_i X_i$$



$$\frac{1}{1 + e^{W^T X}}$$

- Substituting back

$$\operatorname{argmax}_W \sum_l y_c^l (W^T X) + \ln \frac{1}{1 + e^{W^T X}}$$

- Using $\ln(1/a) = -\ln(a)$ \Rightarrow

$$L(W) = \sum_l y_c^l (W^T X) - \ln \left(1 + e^{W^T X} \right)$$

$$\operatorname{argmax}_W L(W) \rightarrow \nabla L(W) = 0$$

$\hookrightarrow \text{solve for } W$

Optimization: Find the max of the log-likelihood function.

- Let's take the gradient of this function w.r.t to W and set it to zero and solve for W.

$$L(W) = \sum_l y_c^l (\underbrace{W^T X}_w) - \ln \left(1 + e^{W^T X} \right)$$

- The gradient is the set of partial derivatives w.r.t each W_i

$$\nabla L(W) = \left(\frac{\partial L(W)}{\partial W_1}, \frac{\partial L(W)}{\partial W_2}, \dots, \frac{\partial L(W)}{\partial W_n} \right)$$

$$\boxed{W^T X = \sum_j w_j x_j}$$

- Set this to zero. Can you solve for W_i ?

$$\frac{\partial L(W)}{\partial W_i} = \sum_l \frac{\partial (y_c^l (\underline{W^T X}))}{\partial W_i} - \frac{\partial}{\partial W_i} \left(\ln \left(1 + e^{W^T X} \right) \right)$$

↓

$$\frac{\partial y_c^l}{\partial W_i} \cdot \frac{\partial \underline{W^T X}}{\partial W_i}$$

$$= \frac{1}{(1 + e^{W^T X})} \cdot \frac{\partial \underline{W^T X}}{\partial W_i}$$

↓

$$= \frac{1}{(1 + e^{W^T X})} \cdot e^{W^T X} \cdot x_i$$

$$\frac{\partial L(W)}{\partial W_i} = \sum_l y_c^l \underbrace{x_i}_{\frac{\partial \underline{W^T X}}{\partial W_i}} - \frac{1}{1 + e^{W^T X}} \cdot e^{W^T X} \cdot x_i$$

$$= 0 \text{ solve for } w_i$$

Optimization: Find the max of the log-likelihood function.

- Let's take the gradient of this function w.r.t to W

$$L(W) = \sum_l y_c^l (W^T X) - \ln \left(1 + e^{W^T X} \right)$$

$$\frac{d}{d\gamma} \ln(\text{prob}) = \frac{1}{\pi} \frac{df(\gamma)}{d\gamma}$$

- This is the set of partial derivatives w.r.t each W_i

$$\nabla L(W) = \left(\frac{\partial L(W)}{\partial W_1}, \frac{\partial L(W)}{\partial W_2}, \dots, \frac{\partial L(W)}{\partial W_n} \right)$$

$$\frac{\partial}{\partial w_i} \ln \left(\frac{1+e^{W^T X}}{1+e^{W^T X}} \right) = \frac{1}{1+e^{W^T X}} \cdot \frac{\partial}{\partial w_i} e^{W^T X}$$

- Where,

$$\frac{\partial L(W)}{\partial W_i} = \sum_l y_c^l X_i - \left(\frac{1}{1 + e^{W^T X}} \right) \left(e^{W^T X} \right) X_i$$

$$= \frac{1}{1+e^{W^T X}} \cdot e^{W^T X} \cdot X_i$$

- You can set this to zero but you cannot solve for W_i because of the non-linearity!

Not to worry! Gradient Descent(Ascent) will save us.

Numerical method for the finding the maximum (and the arg-max).

$$\underset{W}{\operatorname{argmax}} \ L(W)$$

Idea:

Start with $\underline{W} \leftarrow <0, 0, \dots, 0>$

Find an update to \underline{W} that will increase $L(\underline{W})$

Iterate until some stopping criterion is met.

Under some conditions this procedure will find the maximum

Intuition:

Gradient of a function at any point = slope of tangent plane at the point

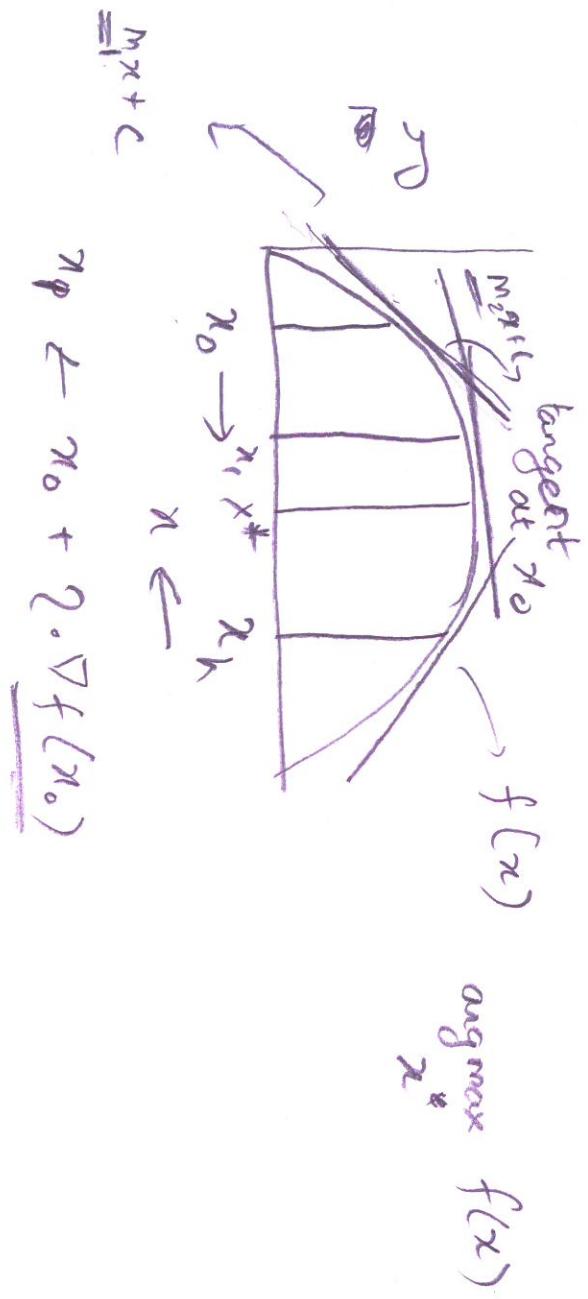
The slope provides the direction in which the function increases most
Take a small step in the gradient direction.

Not to worry! Gradient Descent (Ascent) will save us.

Gradient Ascent Intuition:

Gradient of a function at any point = slope of tangent plane at the point

- The slope provides the direction in which the function increases the most.
- Take a small step in the gradient direction.



Gradient Ascent: Update rule

- Rule

$$W^{(t+1)} \leftarrow W^{(t)} + \eta \cdot \nabla L(W^{(t)})$$

$W^{(t)} \top$

- Recall the partial derivative is:

$$\frac{\partial L(W)}{\partial w_i} = \sum_{\ell} y_{\ell}^l x_i - \boxed{\frac{1}{1+e^{w_i^T x}} \cdot e^{w_i^T x}}$$

Q

$$h_W(x) \triangleq P(y=1|x)$$

- Two things to note:

- Current estimate

$$\frac{\partial L(W)}{\partial w_i} = \sum_{\ell} y_{\ell}^l x_i - \underbrace{P(y=1|x)}_{\substack{\hookrightarrow \text{curr. estimate} \\ \text{for } W^{(t)}}} \cdot x_i$$

- Error term

$$= \sum_{\ell} x_i \left(\underbrace{y_{\ell}^l}_{\substack{\text{true label}}} - \underbrace{P(y=1|x)}_{\substack{\text{curr. estimate} \\ \text{for } W^{(t)}}} \right)$$

True label.

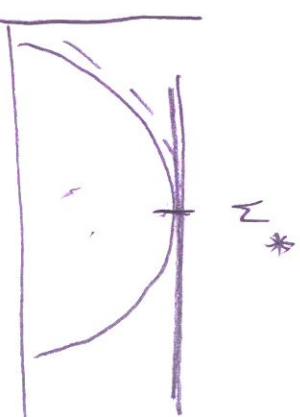
- Update rule per w_i :

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_{\ell} x_i^l \left(\underbrace{y_{\ell}^l}_{\substack{\text{true label}}} - \underbrace{P(y=1|x_i; W^{(t)})}_{\substack{\text{curr. estimate} \\ \text{for } W^{(t)}}} \right)$$

Gradient Ascent: Termination, Convergence, and Alternatives

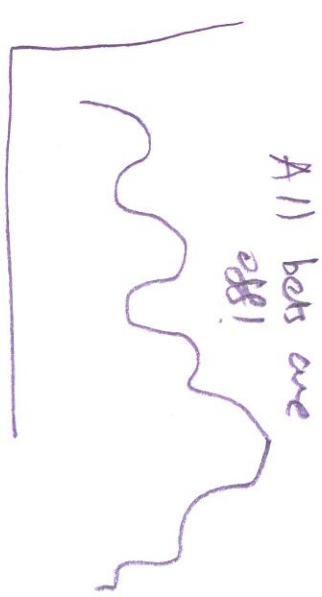
1) When to stop updating?

- > Stop when w doesn't change $\rightarrow w$ will converge b/c at w^* , $w = w^*$
- >



2) Convergence Guarantee

- > Does it converge to the best w ?
- > If $L(w)$ is concave, you will converge to global maximum.



3) Alternatives

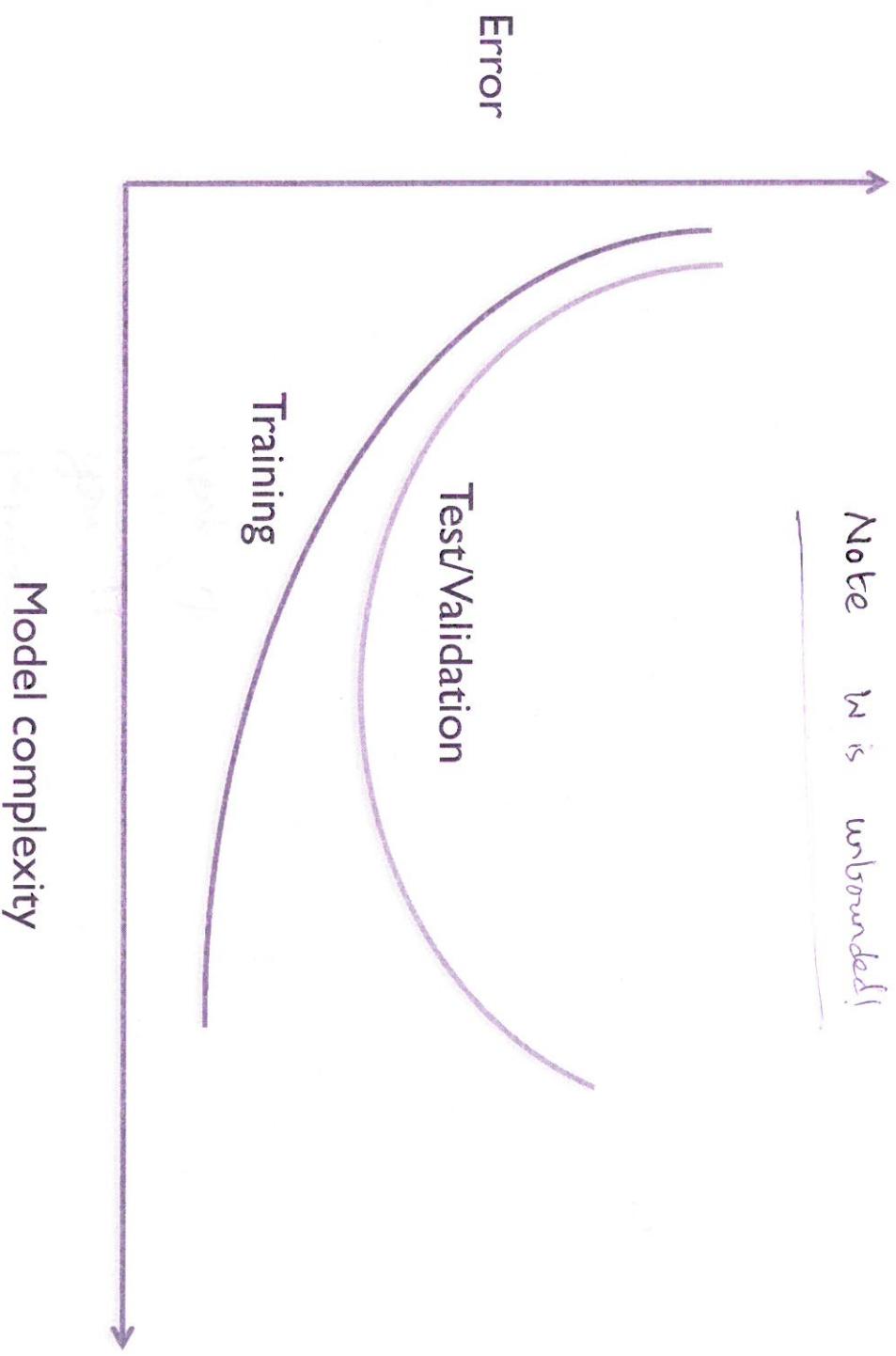
- > Gradient-based optimization methods
BFGS, LBFGS
- > Provide the gradient form & these methods give you the w^*

Regularization

Over-fitting

Training error decreases but not test error!

Note w is unbounded!

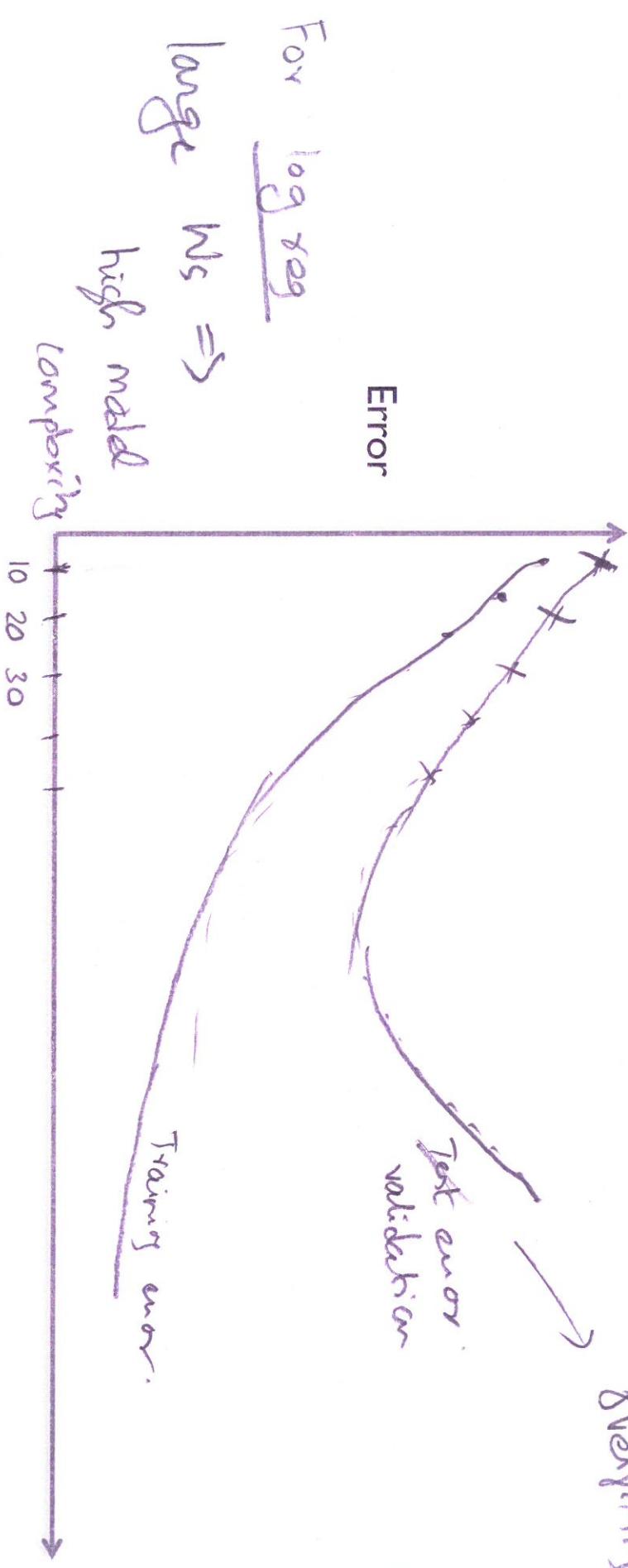


Overfitting and Regularization

Over-fitting: As you add more features i.e., parameters to your model

- The complexity of the model increases.
- The error decreases in training because you have more freedom to fit function.
- At the same time you might be over fitting to some noise
- This means your performance on unseen examples will suffer.

Overfitting!



Regularization

> Issue: Over fitting results in large W values.

> Fix? To insist that W be small.

$$\text{(ie)} \quad \|W\| \rightarrow 0$$

$$\text{argmax}_W L(W) - \frac{\lambda}{2} \|W\|^2$$

> Update rule with regularization

Take gradient

$$\stackrel{\text{alt L1 reg}}{=} \underline{w_i^{(t+1)}}$$

$$\|W\| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$$

$$w_i^2$$

$$2w_i$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum x_i^l (y_c^l - p_r(y=1|x_i^l; W^{(t)}))$$

\uparrow
hyper parameters

Regularization

Over fitting results in large W values.

Regularization aims to enforce sparsity. i.e., small W

Add a penalty term to the Log-likelihood function

$$\underset{W}{\operatorname{argmax}} \ L(W) - \frac{\lambda}{2} \|W\|^2$$

Update rule with regularization parameter

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

e.g. $f_{32}(c, x)$ → is a feature function:
 $f_c(c, x) \leftarrow$ is a feature
 added by c .
 in x that is
 that is present
 in x that is
 others.
 $\cup c = \text{feat.}$

$$\frac{\exp\left(\sum w_i f_i(c, x)\right)}{\exp\left(\sum w_i f_i(c, x)\right)} = P(y=c|x)$$

Multiple classes: $C = \{c_1, c_2, \dots, c_m\}$
 $(x)^M y - 1 =$
 $\wedge P(y=0|x)$
 $(x|y=1) \rightarrow P(x)^M y \approx$
 $P(y|x) \approx h^M(x)$
 Extending to multiple classes

Logistic Regression Summary

- Discriminative version of Naïve Bayes
 - Models the class conditional directly.
 - Both are linear models that are identical when $\overline{w} \rightarrow \infty$
 - Often the first choice when the feature set is not too large.
 - Especially for binary classification.
 - Easily extends to multiple classes
 - Maximum entropy classifier
 - $\log reg \gg NB$ when training data is larger.
 - Useful references
 - Generative and Discriminative Classifiers (Chapter 1), Machine Learning, Tom Mitchell.
 - On Discriminative vs. Generative Classifiers: A comparison of Logistic Regression and Naïve Bayes, Ng and Jordan, NIPS 2002.
 - Andrew Ng's Coursera Lectures on Logistic Regression
 - <https://class.coursera.org/ml-003/lecture>