CSE 573 : Artificial Intelligence

Hanna Hajishirzi Machine Learning, Perceptrons

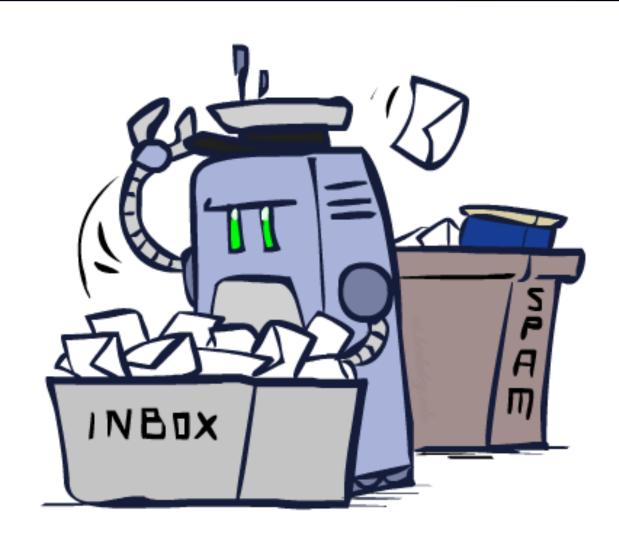
slides adapted from Dan Klein, Pieter Abbeel ai.berkeley.edu And Dan Weld, Luke Zettlemoyer



Machine Learning

- Up until now: how use a model to make optimal decisions
- Machine learning: how to acquire a model from data / experience
 - Learning parameters (e.g. probabilities)
 - Learning structure (e.g. graphs)
 - Learning hidden concepts (e.g. clustering)
- First: model-based classification

Classification



Example: Spam Filter

- Input: an email
- Output: spam/ham
- Setup:
 - Get a large collection of example emails, each labeled "spam" or "ham"
 - Note: someone has to hand label all this data!
 - Want to learn to predict labels of new, future emails
- Features: The attributes used to make the ham / spam decision
 - Words: FREE!
 - Text Patterns: \$dd, CAPS
 - Non-text: SenderInContacts, WidelyBroadcast
 - ...



Dear Sir.

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



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

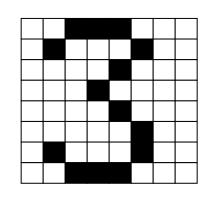
99 MILLION EMAIL ADDRESSES FOR ONLY \$99



Ok, Iknow 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.

Example: Digit Recognition

- Input: images / pixel grids
- Output: a digit 0-9



- 0

- Setup:
 - Get a large collection of example images, each labeled with a digit
 - Note: someone has to hand label all this data!
 - Want to learn to predict labels of new, future digit images

- L
- / 1

- Features: The attributes used to make the digit decision
 - Pixels: (6,8)=ON
 - Shape Patterns: NumComponents, AspectRatio, NumLoops
 - **...**



Other Classification Tasks

Classification: given inputs x, predict labels (classes) y

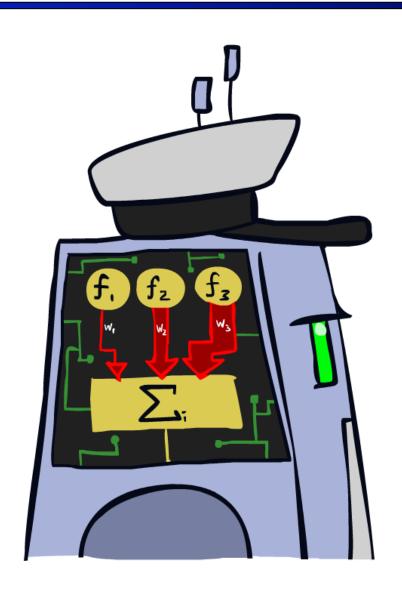
• Examples:

- Spam detection input: document; classes: spam / ham
- OCR input: images; classes: characters
- Medical diagnosis input: symptoms; classes: diseases
- Automatic essay grading input: document; classes: grades
- Fraud detection input: account activity; classes: fraud / no fraud
- Customer service email routing
- ... many more



Classification is an important commercial technology!

Linear Classifiers



A Spam Filter

Data:

- Collection of emails, labeled spam or ham
- Note: someone has to hand label all this data!
- Split into training, heldout, test sets

Classifiers

- Learn on the training set
- (Tune it on a held-out set)
- Test it on new emails



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virture of its nature as being utterly confidencial 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



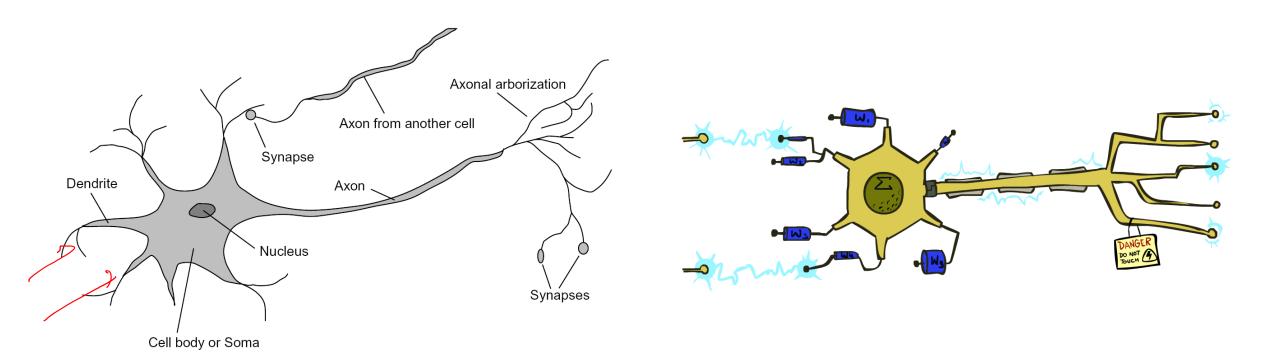
Ok, Iknow this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

Feature Vectors

f(x)# free : 2
YOUR_NAME : 0
MISSPELLED : 2 Hello, **SPAM** Do you want free printr or cartriges? Why pay more when you can get them ABSOLUTELY FREE! Just PIXEL-7,12 : 1
PIXEL-7,13 : 0
...
NUM_LOOPS : 1
...

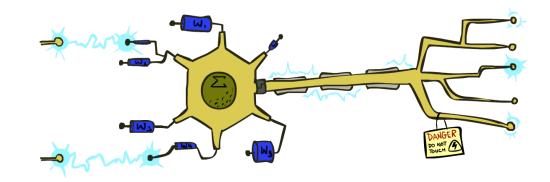
Some (Simplified) Biology

Very loose inspiration: human neurons



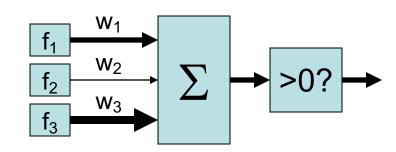
Linear Classifiers

- Inputs are feature values
- Each feature has a weight
- Sum is the activation



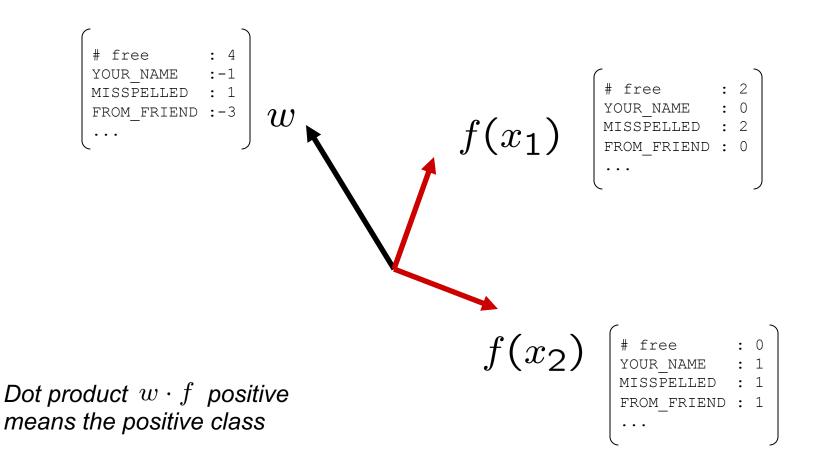
$$activation_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
 - Positive, output +1
 - Negative, output -1

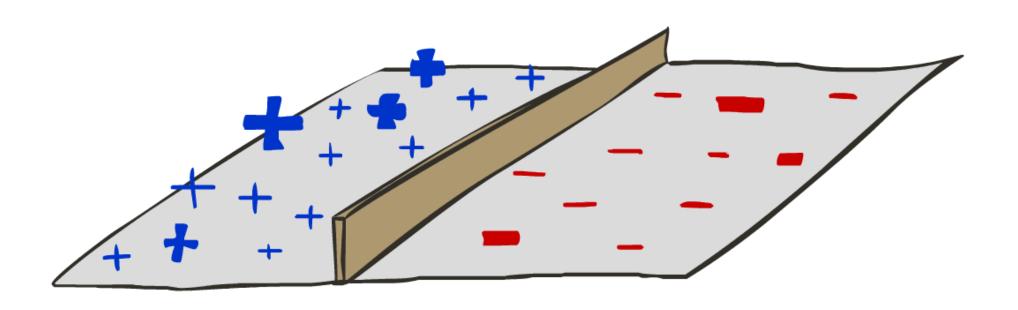


Weights

- Binary case: compare features to a weight vector
- Learning: figure out the weight vector from examples



Decision Rules

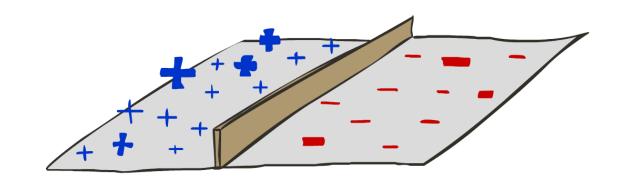


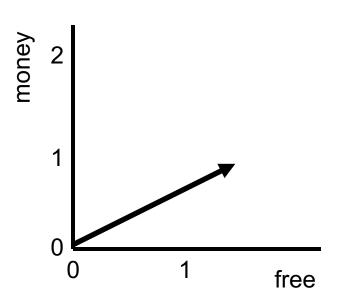
Binary Decision Rule

- In the space of feature vectors
 - Examples are points
 - Any weight vector is a hyperplane
 - One side corresponds to Y=+1
 - Other corresponds to Y=-1

w

BIAS : -3
free : 4
money : 2



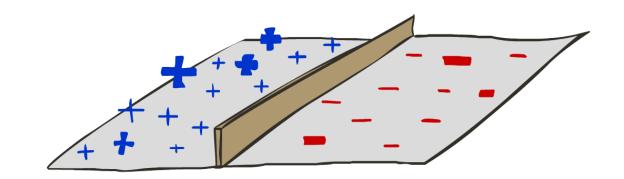


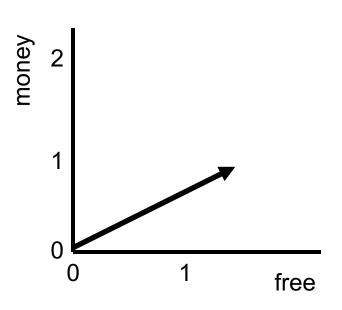
Binary Decision Rule

- In the space of feature vectors
 - Examples are points
 - Any weight vector is a hyperplane
 - One side corresponds to Y=+1
 - Other corresponds to Y=-1

w

free : 4 money : 2



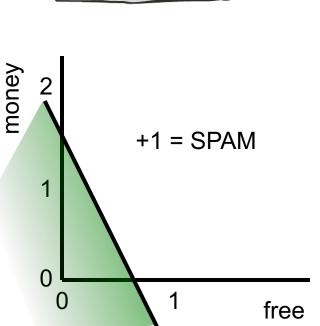


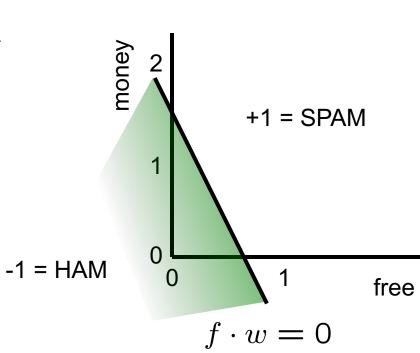
Binary Decision Rule

- In the space of feature vectors
 - Examples are points
 - Any weight vector is a hyperplane
 - One side corresponds to Y=+1
 - Other corresponds to Y=-1

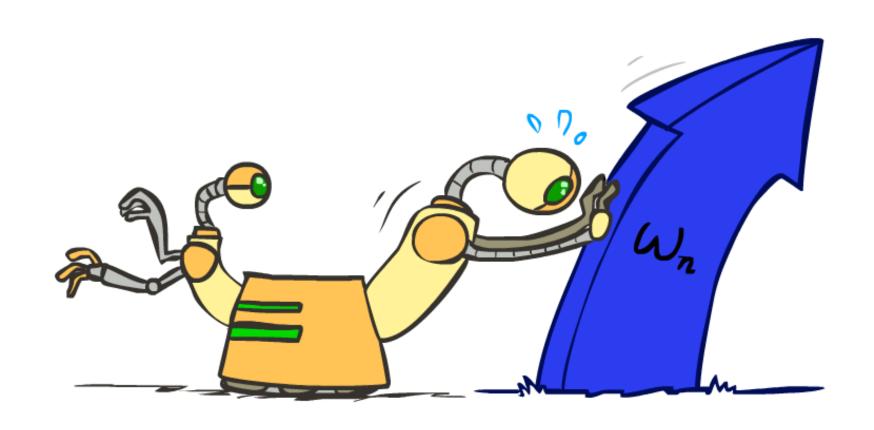
w

BIAS free money:

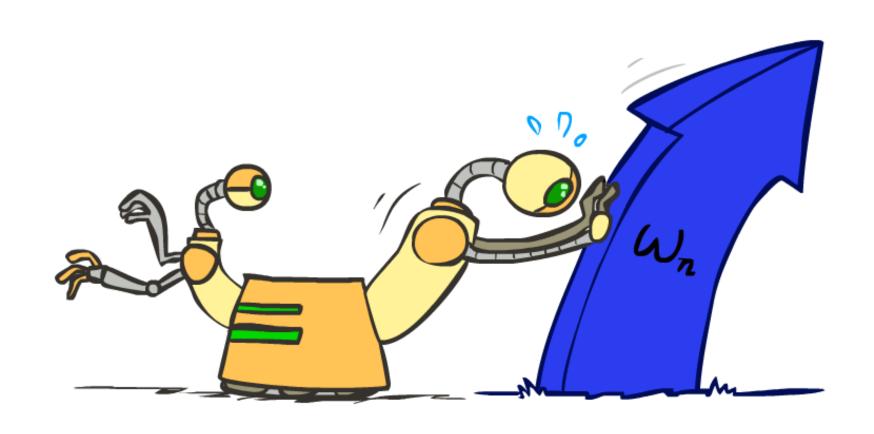




How do we learn weight?



Weight Updates

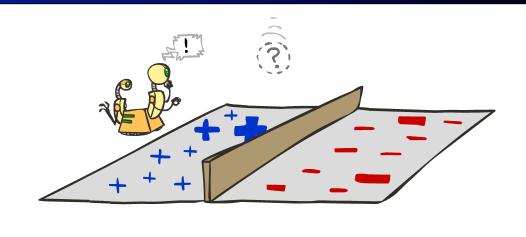


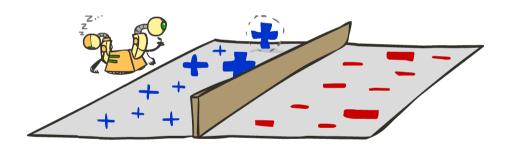
Learning: Binary Perceptron

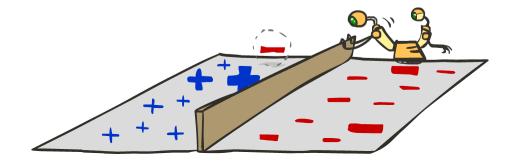
- Start with weights = 0
- For each training instance:
 - Classify with current weights

■ If correct (i.e., y=y*), no change!

• If wrong: adjust the weight vector







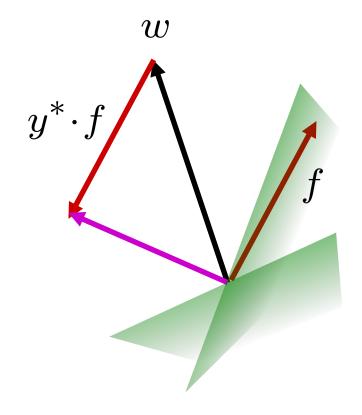
Learning: Binary Perceptron

- Start with weights = 0
- For each training instance:
 - Classify with current weights

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \ge 0\\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$

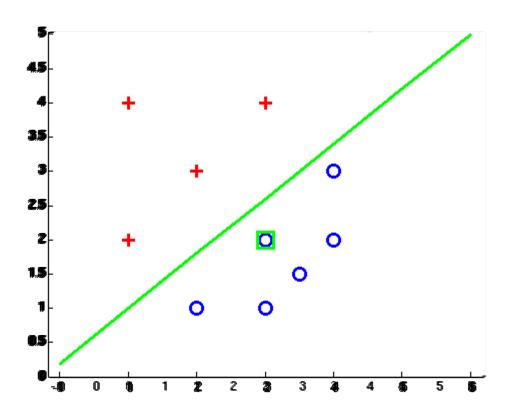
- If correct (i.e., y=y*), no change!
- If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if y* is -1.

$$w = w + y^* \cdot f$$



Examples: Perceptron

Separable Case



Multiclass Decision Rule

- If we have multiple classes:
 - A weight vector for each class:

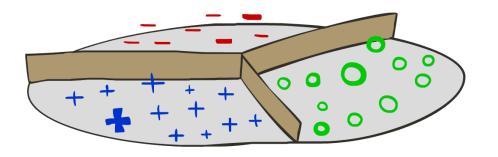
$$w_y$$

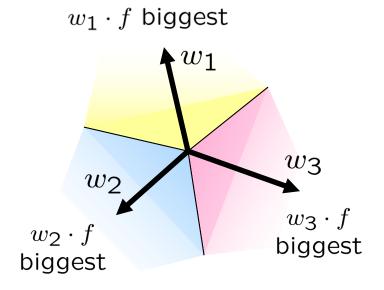
Score (activation) of a class y:

$$w_y \cdot f(x)$$

Prediction highest score wins

$$y = \underset{y}{\operatorname{arg\,max}} \ w_y \cdot f(x)$$





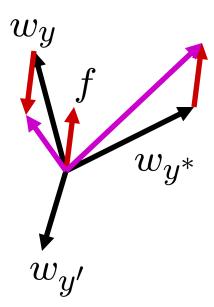
Learning: Multiclass Perceptron

- Start with all weights = 0
- Pick up training examples one by one
- Predict with current weights

$$y = \arg \max_{y} w_{y} \cdot f(x)$$

- If correct, no change!
- If wrong: lower score of wrong answer, raise score of right answer

$$w_y = w_y - f(x)$$
$$w_{y^*} = w_{y^*} + f(x)$$



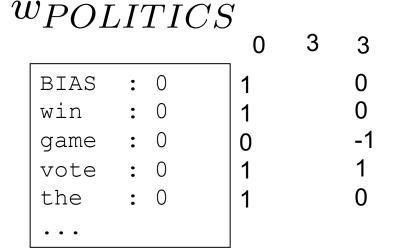
Example: Multiclass Perceptron

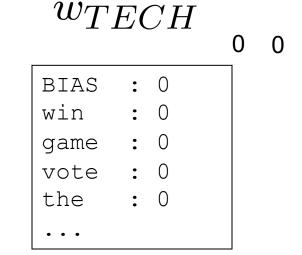
```
"win the vote" [1 1 0 1 1]
```

"win the election" [1 1 0 0 1]

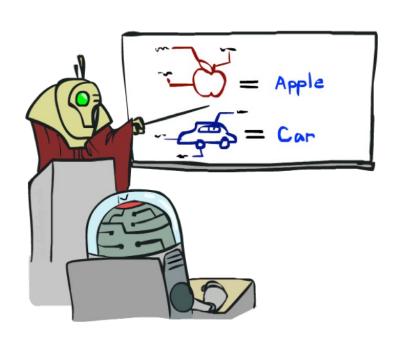
"win the game" [1 1 1 0 1]

w_{SPORTS} 1 -2 -2 BIAS : 1 0 1 win : 0 -1 0 game : 0 0 1 vote : 0 -1 -1 the : 0 -1 0

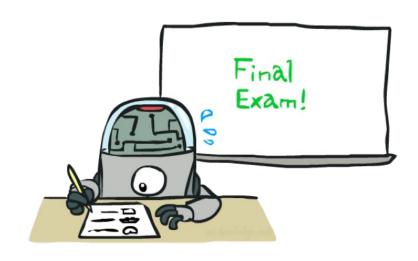




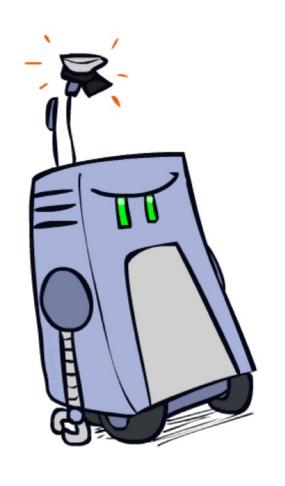
Training and Testing

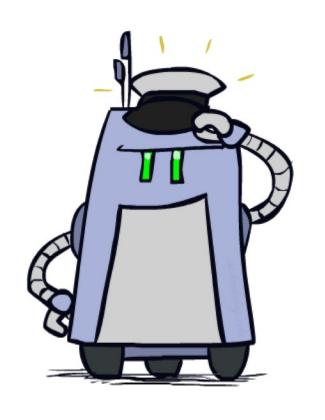


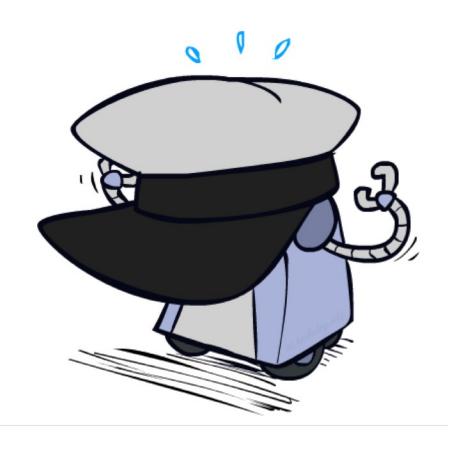




Underfitting and Overfitting







Overfitting

Too many features

- Spam if contains "FREE!"
- Spam if contains \$dd, CAPS
- ...
- Spam if contains "Sir"
- Spam if contains address
- Spam if contains "OT"
- •



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virture of its nature as being utterly confidencial 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, Iknow 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.

Example: Overfitting

$$P(\text{features}, C = 2)$$

$$P(C = 2) = 0.1$$

$$P(\text{on}|C=2) = 0.8$$

P(on|C=2) = 0.1

P(off|C=2) = 0.1

 $P(\mathsf{on}|C=2) = 0.01$



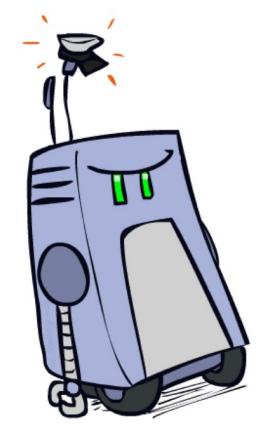
$$P(C = 3) = 0.1$$

$$P(\text{on}|C=3)=0.8$$

$$-P(\text{on}|C=3)=0.9$$

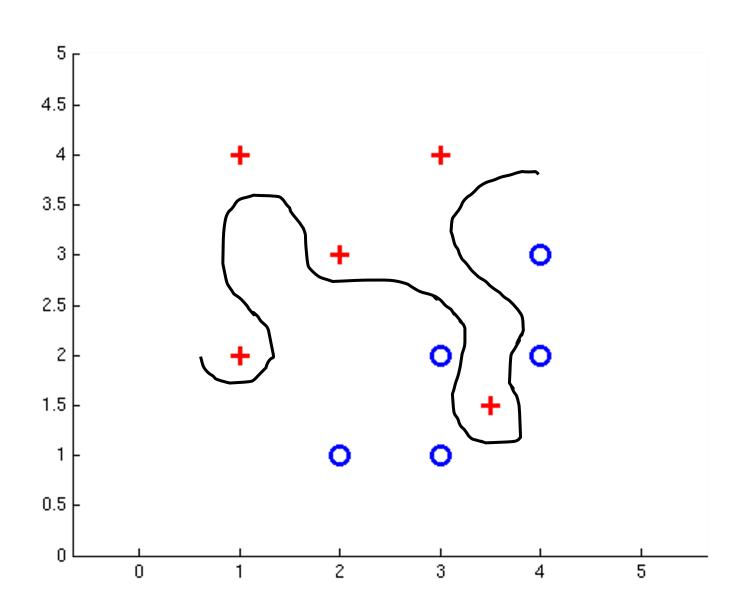
$$P(\text{off}|C=3) = 0.7$$

$$-P(\text{on}|C=3)=0.0$$

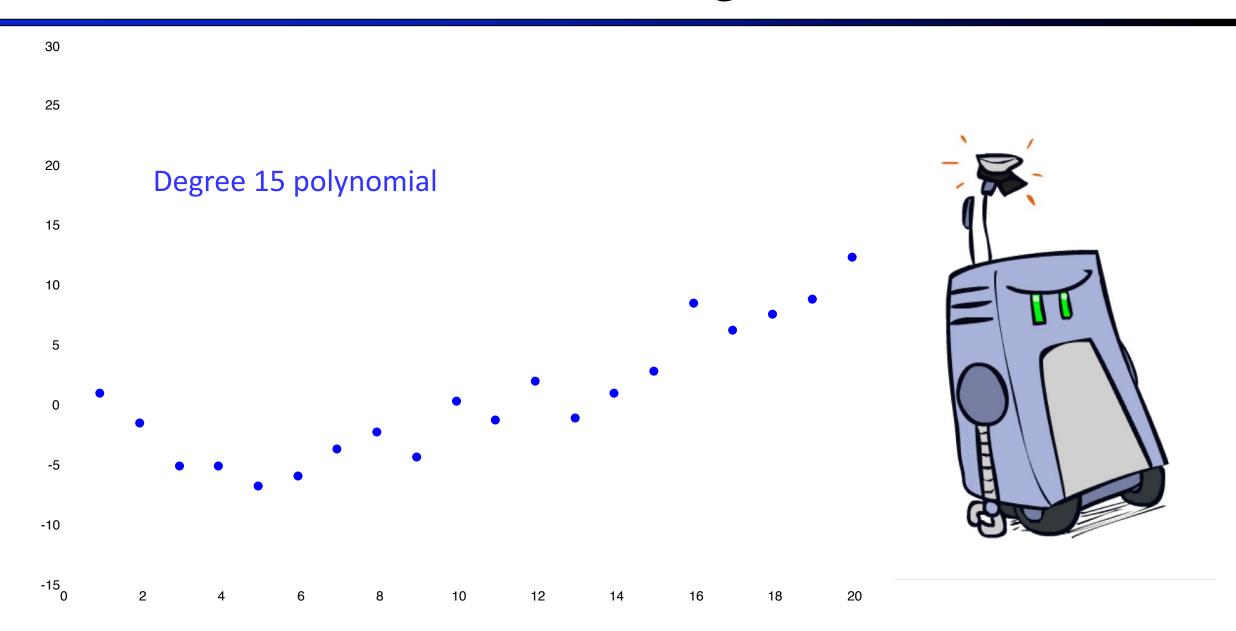


2 wins!!

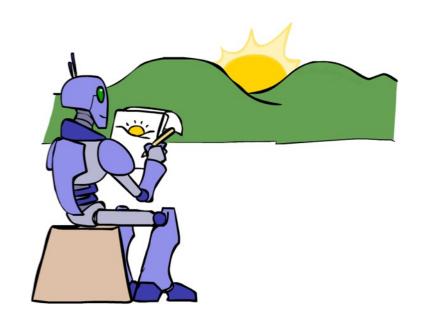
Overfitting

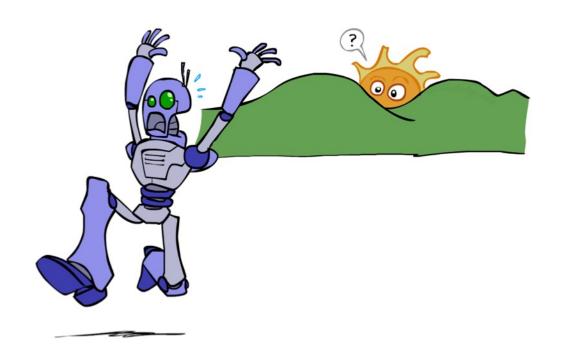


Overfitting



Unseen Events





Generalization and Overfitting

- Relative frequency parameters will overfit the training data!
 - Just because we never saw a non-spam email with an address during training doesn't mean we won't see it at test time
 - 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

Regularization

- Limit the number of features
- Limit the norm of the vectorw
 - If $\mathbf{w_1}$ and $\mathbf{w_2}$ are equally good, and $|\mathbf{w_1}| > |\mathbf{w_2}|$, then $\mathbf{w_2}$ is likely to better generalize

```
FROM_FRIEND :-6
# free : 4
YOUR NAME :-1
MISSPELLED: 1
FROM_FRIEND:-3
```

 $y = \operatorname{arg\,max}_y w_y \cdot f(x) - |w_y|$

$$f(x_2)$$
 $\left(egin{array}{ccccc} \# & { t free} & : & 0 \ { t YOUR_NAME} & : & 1 \ { t MISSPELLED} & : & 1 \ { t FROM_FRIEND} & : & 1 \ { t \cdots} \end{array}
ight)$

Important Concepts

- Data: labeled instances, e.g. emails marked spam/ham
 - Training set
 - Held out set
 - Test set
- Features: attribute-value pairs which characterize each x
- Experimentation cycle
 - Learn parameters (e.g. model probabilities) on training set
 - (Tune hyperparameters on held-out set)
 - Compute accuracy on test set
 - Very important: never "peek" at the test set!
- Evaluation
 - Accuracy: fraction of instances predicted correctly
- Overfitting and generalization
 - Want a classifier which does well on test data
 - Overfitting: fitting the training data very closely, but not generalizing well
 - <u>Underfitting</u>: fits the training set poorly

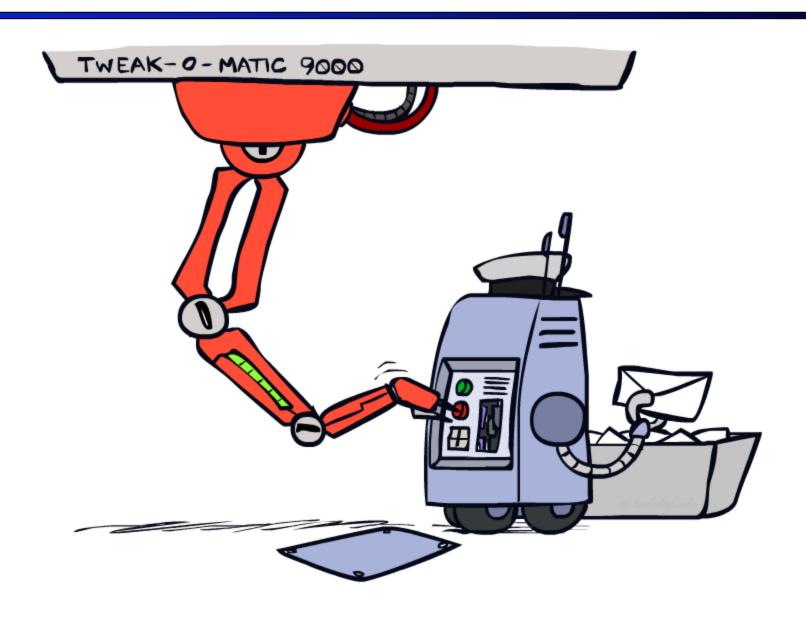
Training
Data

Held-Out Data

> Test Data

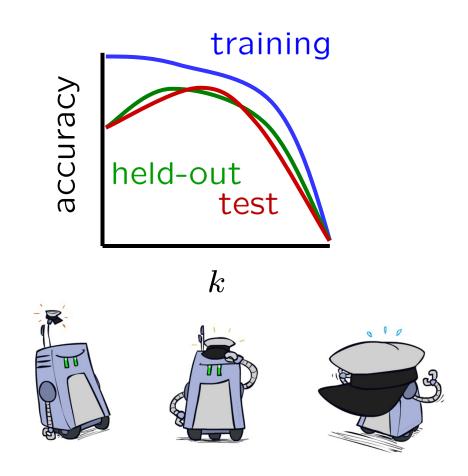


Tuning



Tuning on Held-Out Data

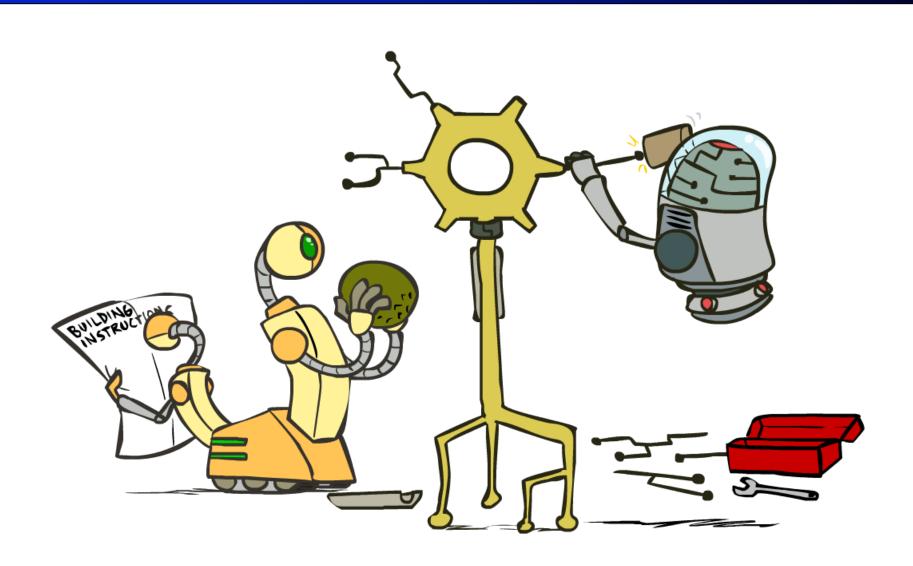
- Now we've got two kinds of unknowns
 - Parameters: the probabilities P(X | Y), P(Y)
 - Hyperparameters: e.g. the amount / type of smoothing to do, k, α
- What should we learn where?
 - Learn parameters from training data
 - Tune hyperparameters on different data
 - Why?
 - For each value of the hyperparameters, train and test on the held-out data
 - Choose the best value and do a final test on the test data



Practical Tip: Baselines

- First step: get a baseline
 - Baselines are very simple "straw man" procedures
 - Help determine how hard the task is
 - Help know what a "good" accuracy is
- Weak baseline: most frequent label classifier
 - Gives all test instances whatever label was most common in the training set
 - E.g. for spam filtering, might label everything as ham
 - Accuracy might be very high if the problem is skewed
 - E.g. calling everything "ham" gets 66%, so a classifier that gets 70% isn't very good...
- For real research, usually use previous work as a (strong) baseline

Improving the Perceptron

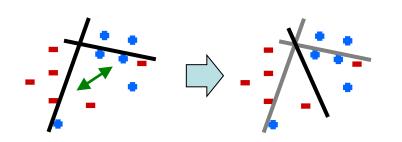


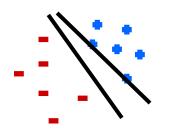
Problems with the Perceptron

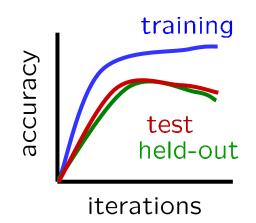
- Noise: if the data isn't separable, weights might thrash
 - Averaging weight vectors over time can help (averaged perceptron)

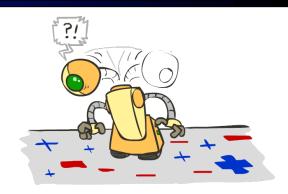


- Overtraining: test / held-out accuracy usually rises, then falls
 - Overtraining is a kind of overfitting

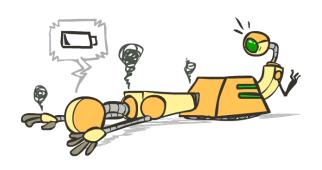




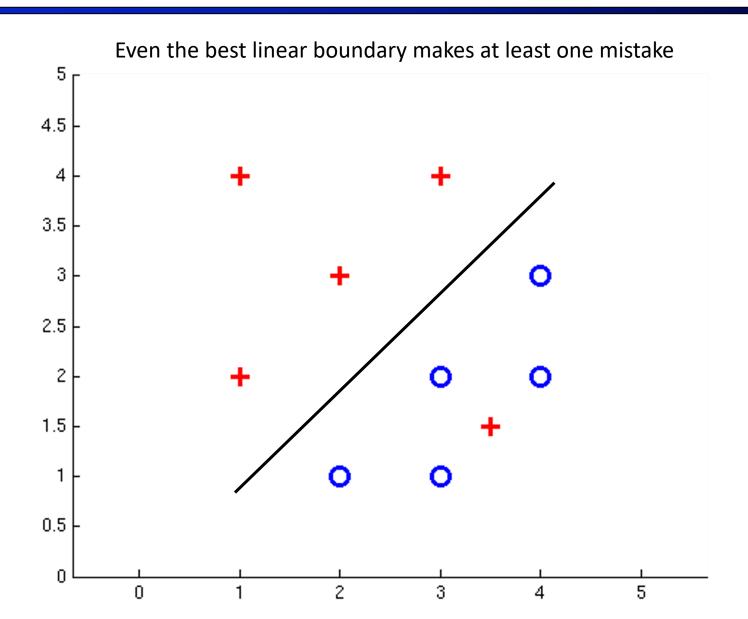




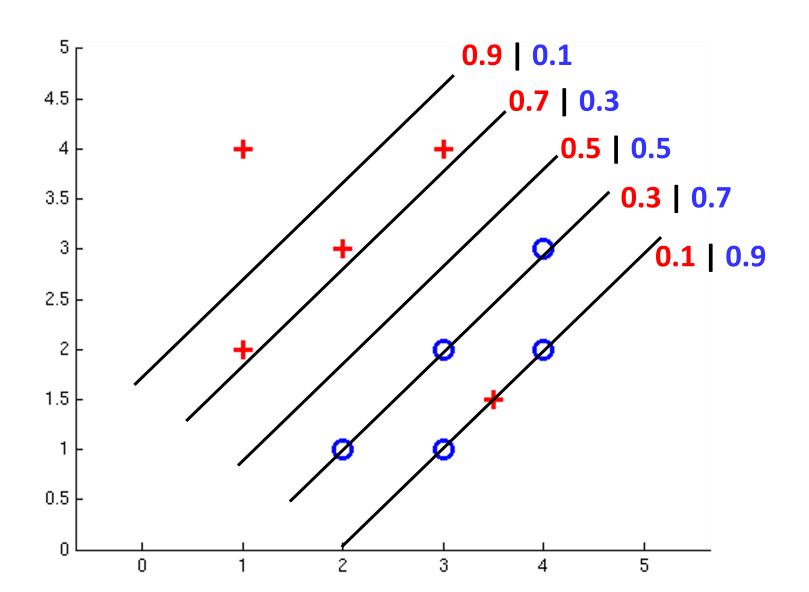




Non-Separable Case: Deterministic Decision



Non-Separable Case: Probabilistic Decision

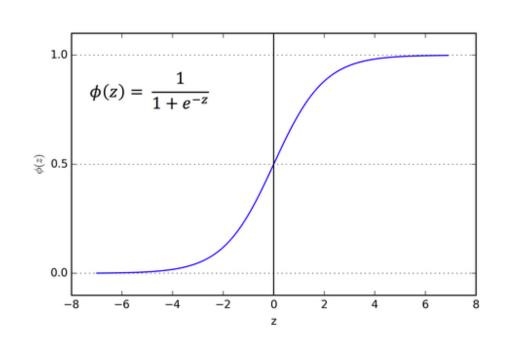


How to get probabilistic decisions?

- Perceptron scoring: $z = w \cdot f(x)$
- If $z = w \cdot f(x)$ very positive \rightarrow want probability going to 1
- If $z = w \cdot f(x)$ very negative \rightarrow want probability going to 0

Sigmoid function

$$\phi(z) = \frac{1}{1 + e^{-z}}$$



Best w?

• Maximum likelihood estimation:

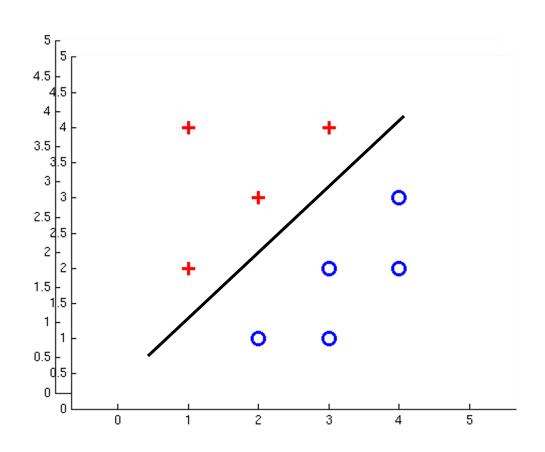
$$\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)}|x^{(i)};w)$$

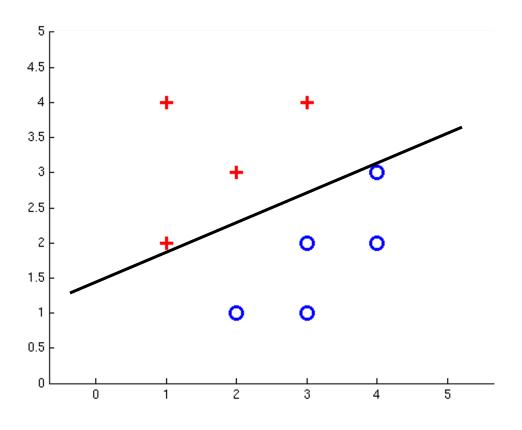
with:
$$P(y^{(i)} = +1|x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

$$P(y^{(i)} = -1|x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

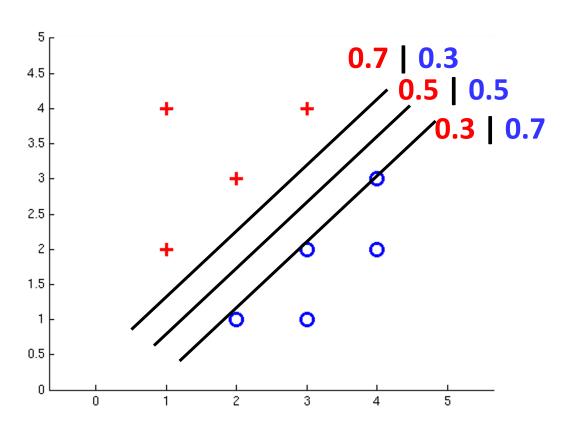
= Logistic Regression

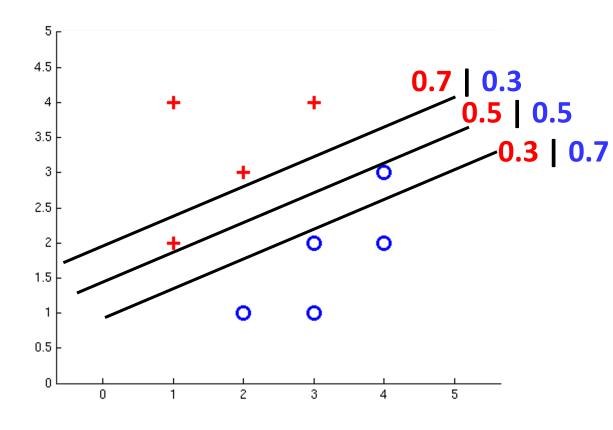
Separable Case: Deterministic Decision – Many Options





Separable Case: Probabilistic Decision – Clear Preference



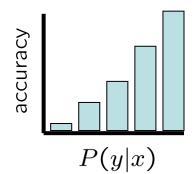


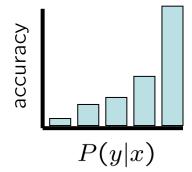
Confidences from a Classifier

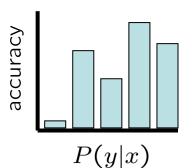
- The confidence of a probabilistic classifier:
 - Posterior over the top label

$$confidence(x) = \max_{y} P(y|x)$$

- Represents how sure the classifier is of the classification
- Any probabilistic model will have confidences
- No guarantee confidence is correct
- Calibration
 - Weak calibration: higher confidences mean higher accuracy
 - Strong calibration: confidence predicts accuracy rate
 - What's the value of calibration?



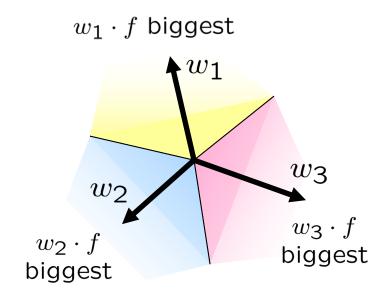




Multiclass Logistic Regression

Recall Perceptron:

- ullet A weight vector for each class: w_y
- Score (activation) of a class y: $w_y \cdot f(x)$
- Prediction highest score wins $y = \arg\max_{y} w_y \cdot f(x)$



• How to make the scores into probabilities?

$$z_1,z_2,z_3 \to \frac{e^{z_1}}{e^{z_1}+e^{z_2}+e^{z_3}}, \frac{e^{z_2}}{e^{z_1}+e^{z_2}+e^{z_3}}, \frac{e^{z_3}}{e^{z_1}+e^{z_2}+e^{z_3}}, \frac{e^{z_3}}{e^{z_1}+e^{z_2}+e^{z_3}}$$
 original activations softmax activations

Best w?

• Maximum likelihood estimation:

$$\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)}|x^{(i)};w)$$

with:
$$P(y^{(i)}|x^{(i)};w) = \frac{e^{w_{y^{(i)}} \cdot f(x^{(i)})}}{\sum_{y} e^{w_{y} \cdot f(x^{(i)})}}$$

= Multi-Class Logistic Regression

Best w?

- Optimization
 - i.e., how do we solve:

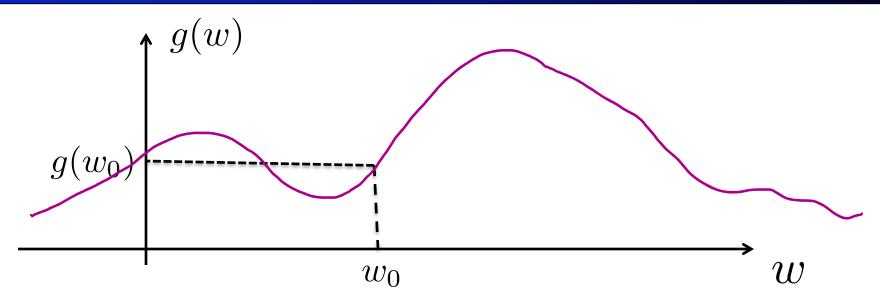
$$\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)}|x^{(i)};w)$$

Hill Climbing

- Simple, general idea
 - Start wherever
 - Repeat: move to the best neighboring state
 - If no neighbors better than current, quit

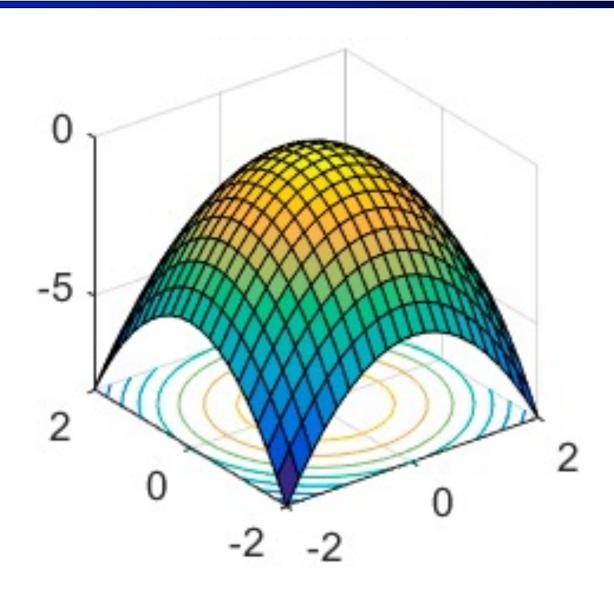
- What's particularly tricky when hill-climbing for multiclass logistic regression?
 - Optimization over a continuous space
 - Infinitely many neighbors!
 - How to do this efficiently?

1-D Optimization



- Could evaluate $g(w_0 + h)$ and $g(w_0 h)$
 - Then step in best direction
- Or, evaluate derivative:
 - Tells which direction to step into

2-D Optimization



Gradient Ascent

- Perform update in uphill direction for each coordinate
- E.g., consider:
 - Updates: $g(w_1, w_2)$

$$w_1 \leftarrow w_1 + \alpha * \frac{\partial g}{\partial w_1}(w_1, w_2)$$

$$w_2 \leftarrow w_2 + \alpha * \frac{\partial g}{\partial w_2}(w_1, w_2)$$

Updates in vector notation:

$$w \leftarrow w + \alpha * \nabla_w g(w)$$

with:
$$\nabla_w g(w) = \begin{bmatrix} \frac{\partial g}{\partial w_1}(w) \\ \frac{\partial g}{\partial w_2}(w) \end{bmatrix}$$
 = gradient

Gradient Ascent

- Idea:
 - Start somewhere
 - Repeat: Take a step in the gradient direction

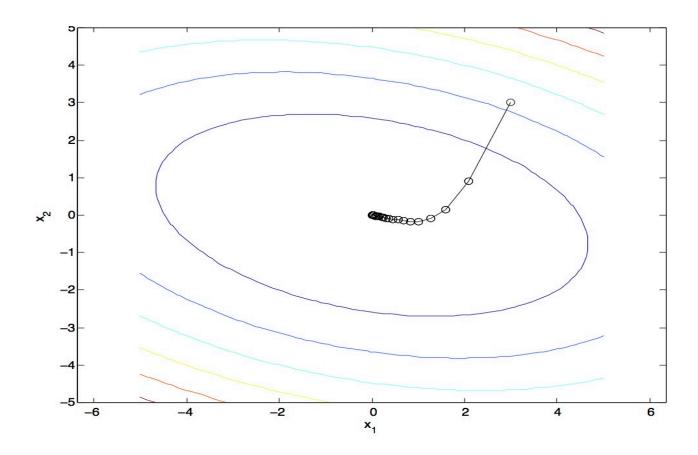


Figure source: Mathworks

Gradient in n dimensions

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \\ \cdots \\ \frac{\partial g}{\partial w_n} \end{bmatrix}$$

Optimization Procedure: Gradient Ascent

• init w• for iter = 1, 2, ... $w \leftarrow w + \alpha * \nabla g(w)$

- ullet α : learning rate --- tweaking parameter that needs to be chosen carefully
- How? Try multiple choices
 - lacktriangle Crude rule of thumb: update changes w about 0.1 1 %

Batch Gradient Ascent on the Log Likelihood Objective

$$\max_{w} ll(w) = \max_{w} \sum_{i} \log P(y^{(i)}|x^{(i)}; w)$$

$$g(w)$$

- lacksquare init w
- for iter = 1, 2, ...

$$w \leftarrow w + \alpha * \sum_{i} \nabla \log P(y^{(i)}|x^{(i)};w)$$

Stochastic Gradient Ascent on the Log Likelihood Objective

$$\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)}|x^{(i)}; w)$$

Observation: once gradient on one training example has been computed, might as well incorporate before computing next one

- lacksquare init w
- for iter = 1, 2, ...
 - pick random j

$$w \leftarrow w + \alpha * \nabla \log P(y^{(j)}|x^{(j)};w)$$

Mini-Batch Gradient Ascent on the Log Likelihood Objective

$$\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)}|x^{(i)}; w)$$

Observation: gradient over small set of training examples (=mini-batch) can be computed in parallel, might as well do that instead of a single one

- lacksquare init w
- for iter = 1, 2, ...
 - pick random subset of training examples J

$$w \leftarrow w + \alpha * \sum_{j \in J} \nabla \log P(y^{(j)} | x^{(j)}; w)$$

How about computing all the derivatives?

 We'll talk about that in neural networks, which are a generalization of logistic regression