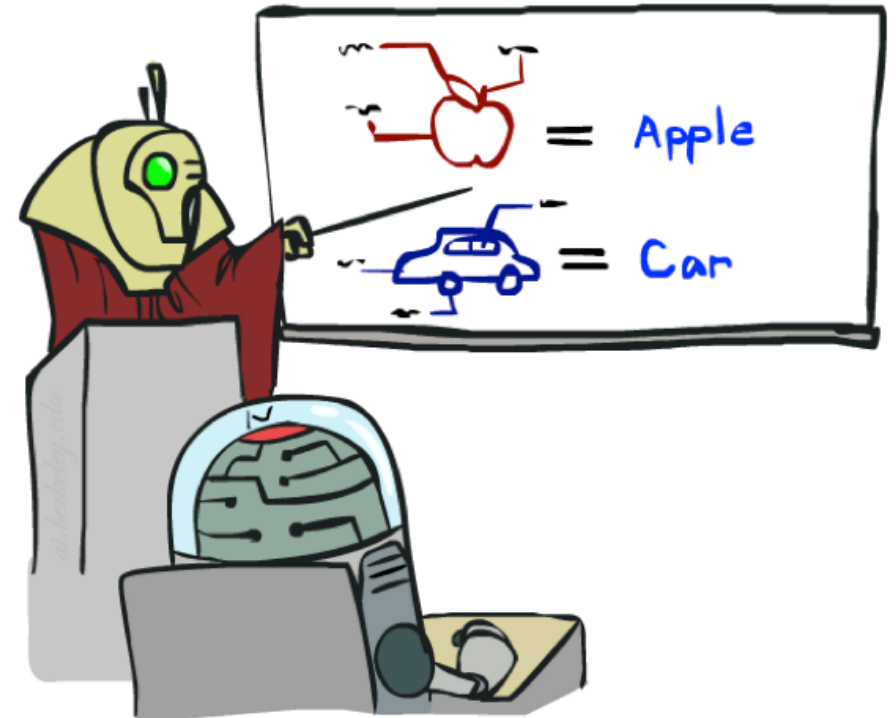# Supervised Machine Learning



AIMA Chapter 18, 20
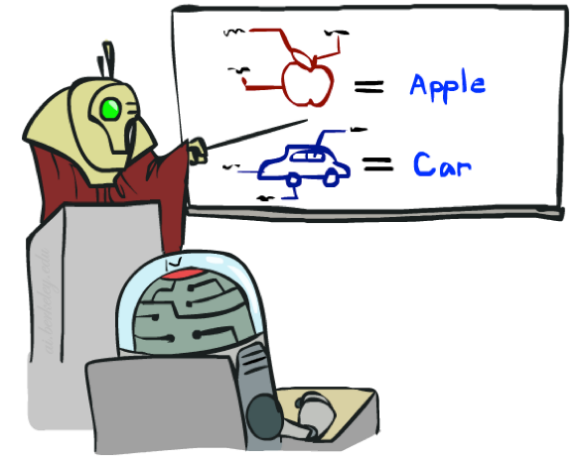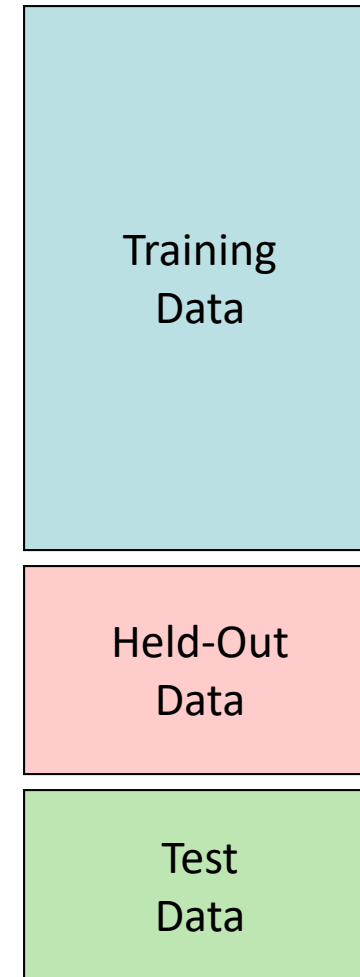
# Supervised learning

- To learn an unknown *target function* f

- Input: a *training set* of *labeled examples* $(x_j, y_j)$ where $y_j = f(x_j)$

- Output: *hypothesis* h that is "close" to f

- Types of supervised learning
  - Classification = learning f with discrete output value
  - Regression = learning f with real-valued output value
  - Structured prediction = learning f with structured output

# Important Concepts

- **Data: labeled instances, e.g. emails marked spam/ham**
  - Training set
  - Held out set
  - Test set

- **Experimentation cycle**
  - Learn parameters (e.g. model probabilities) on training set
  - Tune hyperparameters on held-out set
  - Compute accuracy of test set (fraction of instances predicted correctly)
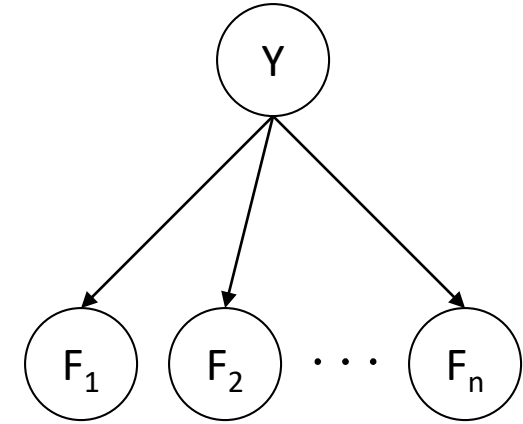  - Very important: never "peek" at the test set!

# Naïve Bayes

- Naïve Bayes model:

|Y| parameters

$$P(\mathsf{Y}, \mathsf{F}_1 \ldots \mathsf{F}_n) = \quad P(\mathsf{Y}) \prod_i P(\mathsf{F}_i | \mathsf{Y})$$
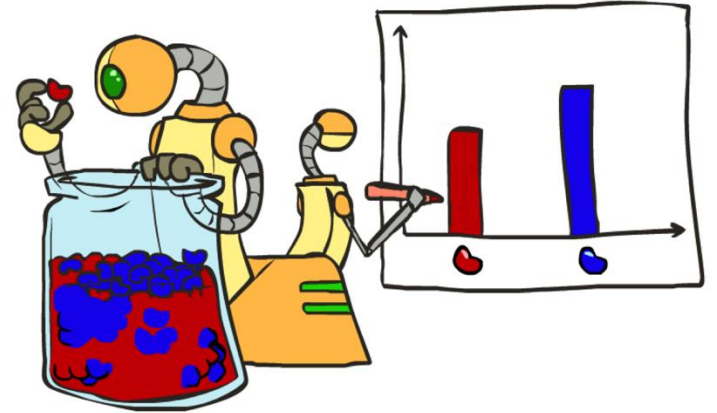
n x |F| x |Y|
parameters



- Assume all features are independent effects of the label

- Total number of parameters is *linear* in n
- Tree-structured: *linear* inference time
- Model is very simplistic, but often works anyway

# Parameter Estimation

- Estimating the distribution of a random variable
- *Elicitation:* ask a human (this is hard...)
- *Empirically:* use training data (learning!)
  - For each outcome x, look at the *empirical rate* of that value

    $$P_{\mathsf{ML}}(x) = \frac{\mathsf{count}(x)}{\mathsf{total\ samples}}$$

  - Ex:
    - We've seen 1000 words from spam emails, among which we see "money" for 50 times
    - So we set P(money | spam) = 0.05

  - This is the estimate that maximizes the *likelihood of the data*
    - Likelihood: conditional probability of the data given the parameters

# Generalization and Overfitting

- **Overfitting: learn to fit the training data very closely, but fit the test data poorly**
  - Generalization: try to fit the test data as well
- **Why does overfitting occur?**
  - Training data is not representative of the true data distribution
    - Too few training samples
    - Training data is noisy
  - Too many attributes, some of them irrelevant to the classification task
  - The model is too expressive
    - Ex: the model is capable of memorizing all the spam emails in the training set

# Generalization and Overfitting

- **Avoid overfitting**
  - Acquire more training data (not always possible)
  - Remove irrelevant attributes (not always possible)
  - Limit the model expressiveness by regularization, early stopping, pruning, etc.

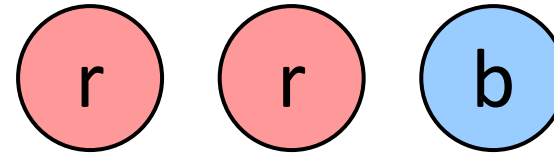- **In our previous example, we may smooth the empirical rate to improve generalization**

# Laplace Smoothing

- Laplace's estimate:
  - Pretend you saw every outcome once more than you actually did

r    r    b

$$P_{LAP}(x) = \frac{c(x) + 1}{\sum_x [c(x) + 1]}$$

$$= \frac{c(x) + 1}{N + |X|}$$

$$P_{ML}(X) =$$

$$P_{LAP}(X) =$$

  - Can derive this estimate with *Dirichlet priors* (see cs281a)

# Laplace Smoothing

- **Laplace's estimate (extended):**
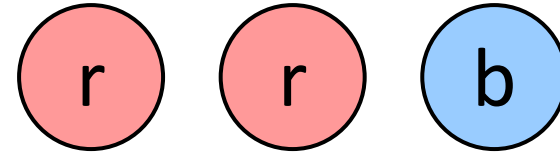  - Pretend you saw every outcome k extra times

$$P_{LAP,k}(x) = \frac{c(x) + k}{N + k|X|}$$

  - What's Laplace with k = 0?
  - k is the <span style="color:red">strength</span> of the prior

- **Laplace for conditionals:**
  - Smooth each condition independently:

$$P_{LAP,k}(x|y) = \frac{c(x,y) + k}{c(y) + k|X|}$$

r r b

$$P_{LAP,0}(X) =$$

$$P_{LAP,1}(X) =$$

$$P_{LAP,100}(X) =$$

# Real NB: Smoothing

- For real classification problems, smoothing is critical
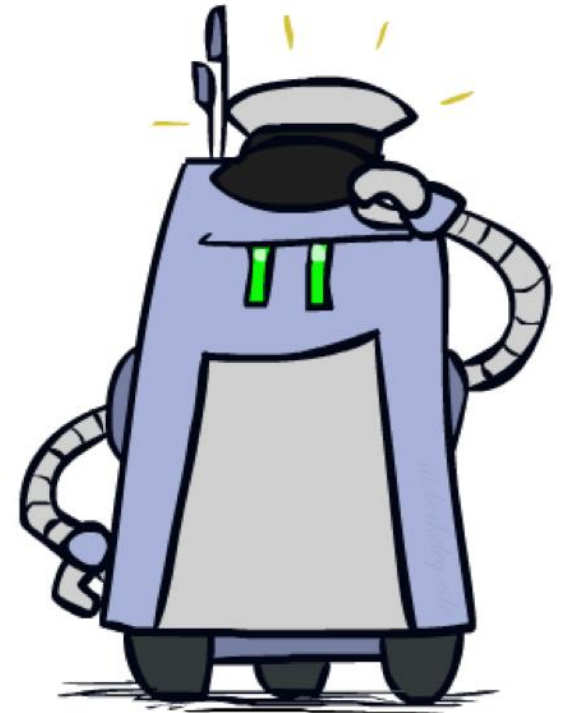
- New odds ratios:

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

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

```
helvetica : 11.4
seems     : 10.8
group     : 10.2
ago       :  8.4
areas     :  8.3
...
```

```
verdana : 28.8
Credit  : 28.4
ORDER   : 27.2
<FONT>  : 26.9
money   : 26.5
...
```
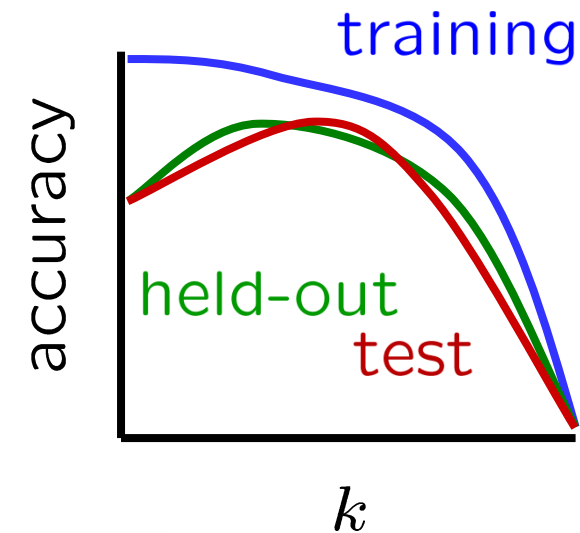
*Do these make more sense?*

# Linear Interpolation

- In practice, Laplace often performs poorly for P(X|Y):
  - When |X| is very large
  - When |Y| is very large

- Another option: linear interpolation
  - Also get the empirical P(X) from the data
  - Make sure the estimate of P(X|Y) isn't too different from the empirical P(X)

$$P_{LIN}(x|y) = \alpha \hat{P}(x|y) + (1.0 - \alpha)\hat{P}(x)$$

# 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, $\alpha$

- **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

# Confidences from a Classifier

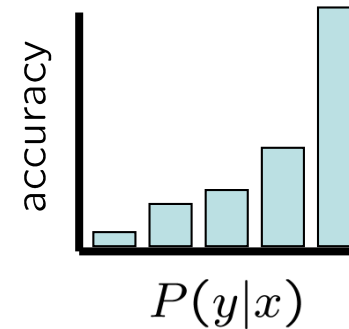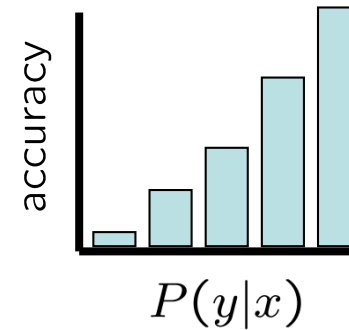- The confidence of a probabilistic classifier:
  - Posterior over the top label

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

  - Represents how sure the classifier is of the classification
  - Any probabilistic model will have confidences
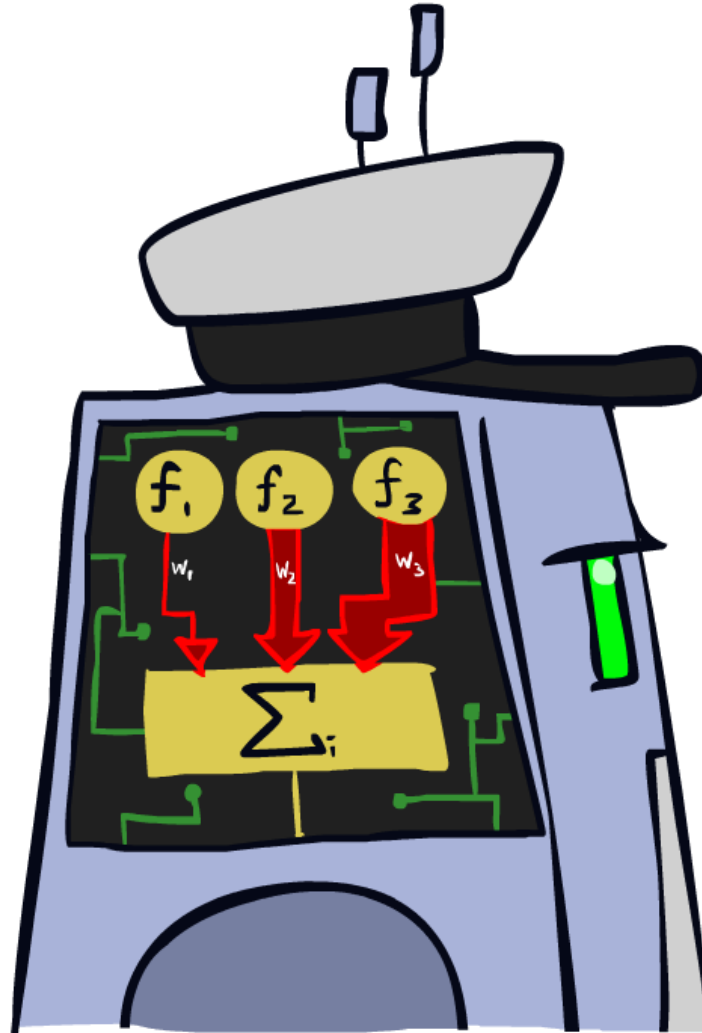  - No guarantee confidence is correct

- Calibration
  - Weak calibration: higher confidences mean higher accuracy
  - Strong calibration: confidence predicts accuracy rate
  - What's the value of calibration?

# Summary

- Bayes rule lets us do diagnostic queries with causal probabilities

- The naïve Bayes assumption takes all features to be independent given the class label

- We can build classifiers out of a naïve Bayes model using training data

- Smoothing estimates is important in real systems

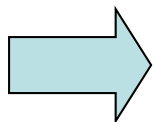- Classifier confidences are useful, when you can get them
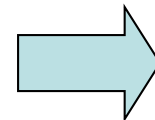
# Linear Classifiers

# Feature Vectors

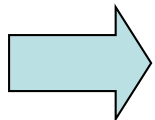$$x \qquad\qquad f(x) \qquad\qquad y$$

```
Hello,

Do you want free printr
cartriges?  Why pay more
when you can get them
ABSOLUTELY FREE!  Just
```

$\Rightarrow$

```
# free      : 2
YOUR_NAME   : 0
MISSPELLED  : 2
FROM_FRIEND : 0
...
```
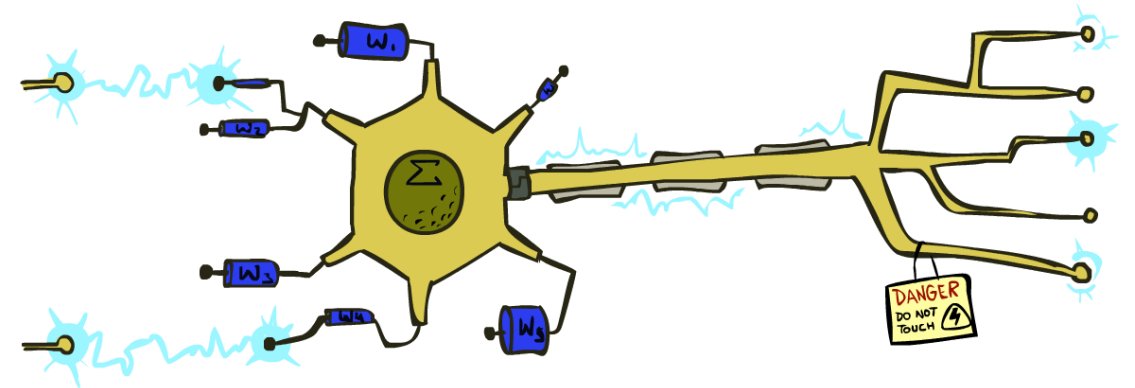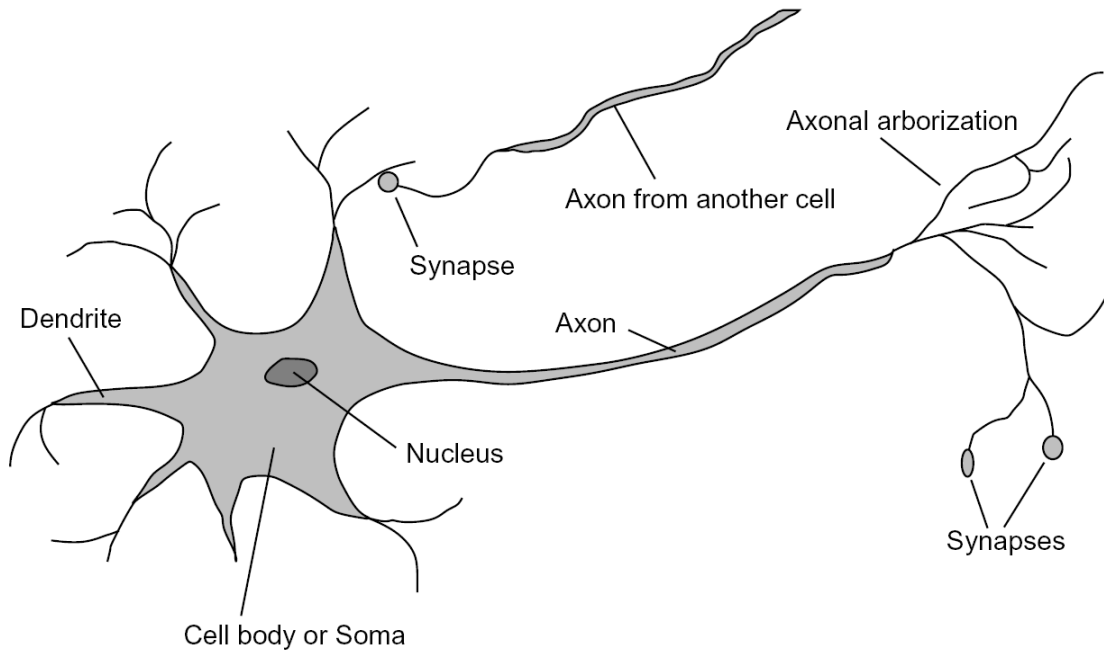
$\Rightarrow$

SPAM
or
+

$\Rightarrow$

```
PIXEL-7,12  : 1
PIXEL-7,13  : 0
...
NUM_LOOPS   : 1
...
```
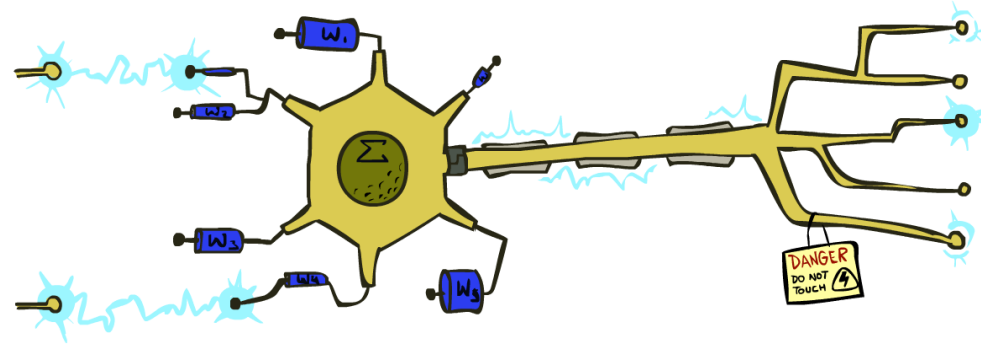
$\Rightarrow$

"2"

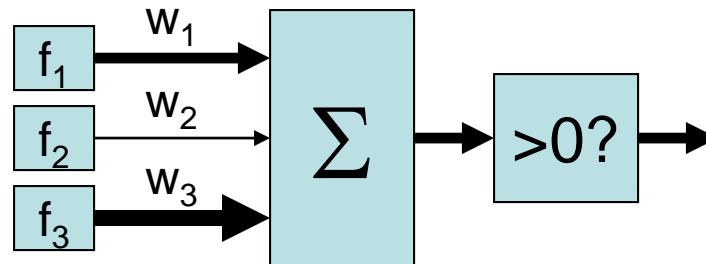# Some (Simplified) Biology

- Very loose inspiration: human neurons

# Linear Classifiers

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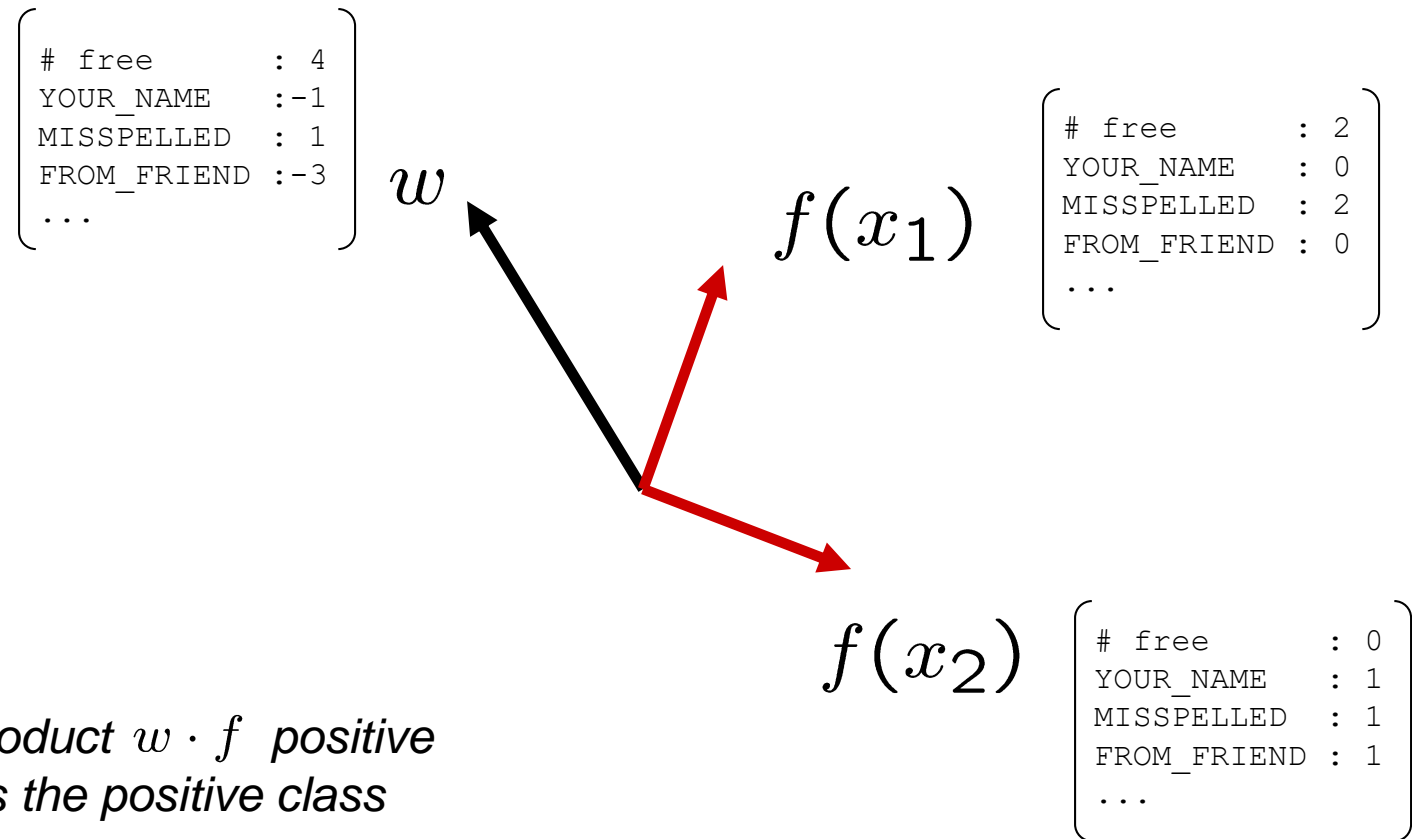- Each feature has a weight
- Sum is the activation



$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- 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

$$\begin{bmatrix} \texttt{\# free : 4} \\ \texttt{YOUR\_NAME :-1} \\ \texttt{MISSPELLED : 1} \\ \texttt{FROM\_FRIEND :-3} \\ \texttt{...} \end{bmatrix} w$$

$$f(x_1) \begin{bmatrix} \texttt{\# free : 2} \\ \texttt{YOUR\_NAME : 0} \\ \texttt{MISSPELLED : 2} \\ \texttt{FROM\_FRIEND : 0} \\ \texttt{...} \end{bmatrix}$$

$$f(x_2) \begin{bmatrix} \texttt{\# free : 0} \\ \texttt{YOUR\_NAME : 1} \\ \texttt{MISSPELLED : 1} \\ \texttt{FROM\_FRIEND : 1} \\ \texttt{...} \end{bmatrix}$$

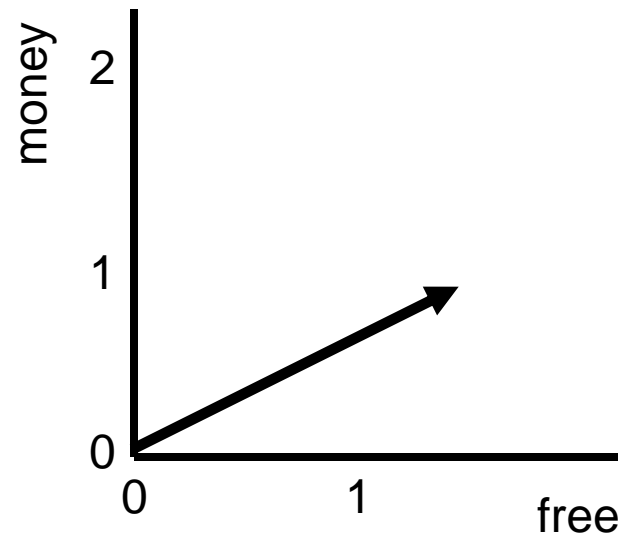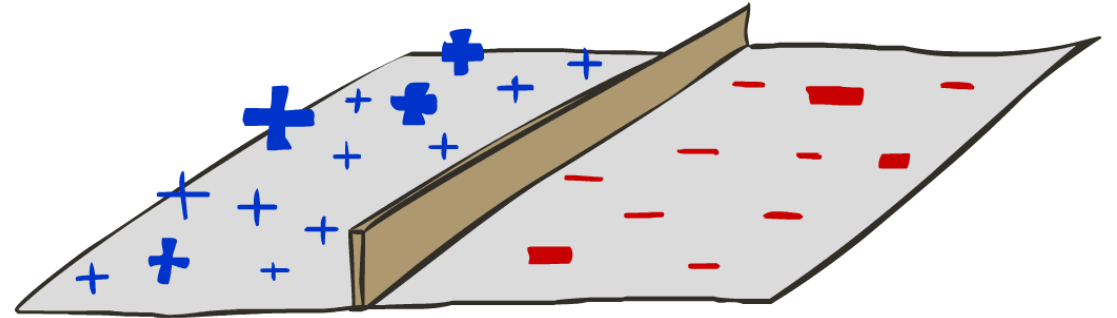*Dot product $w \cdot f$ positive means the positive class*

# 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
...
```
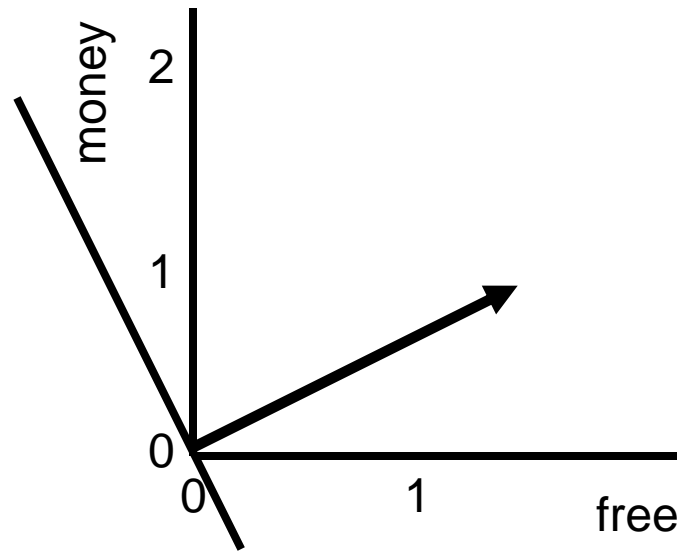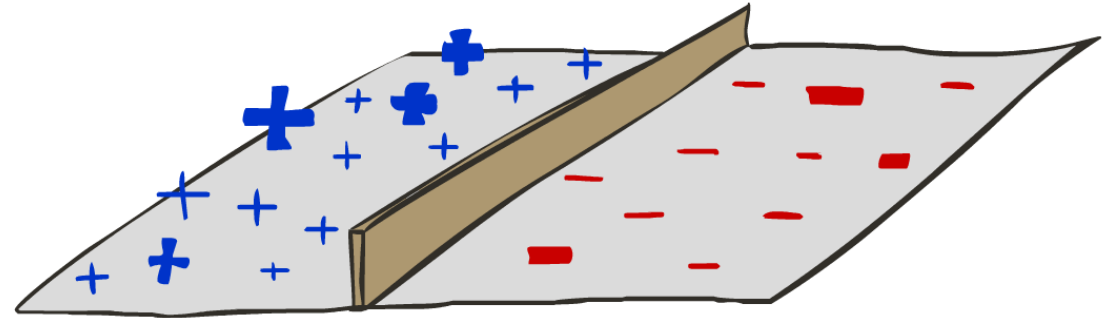
# Binary Decision Rule

- **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

$w$

```
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
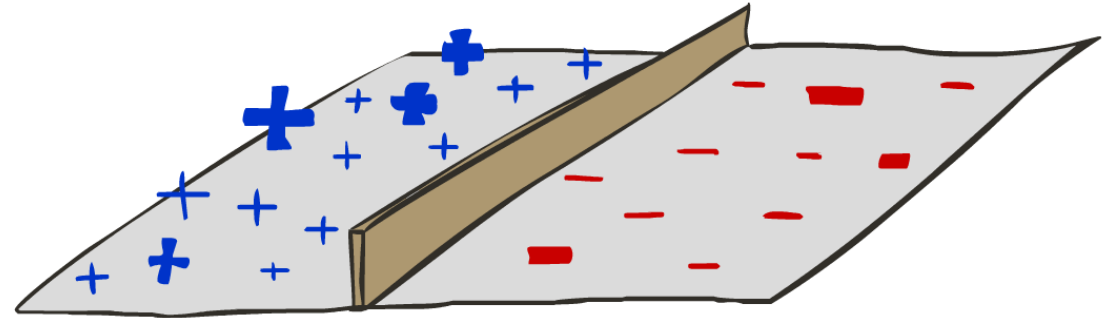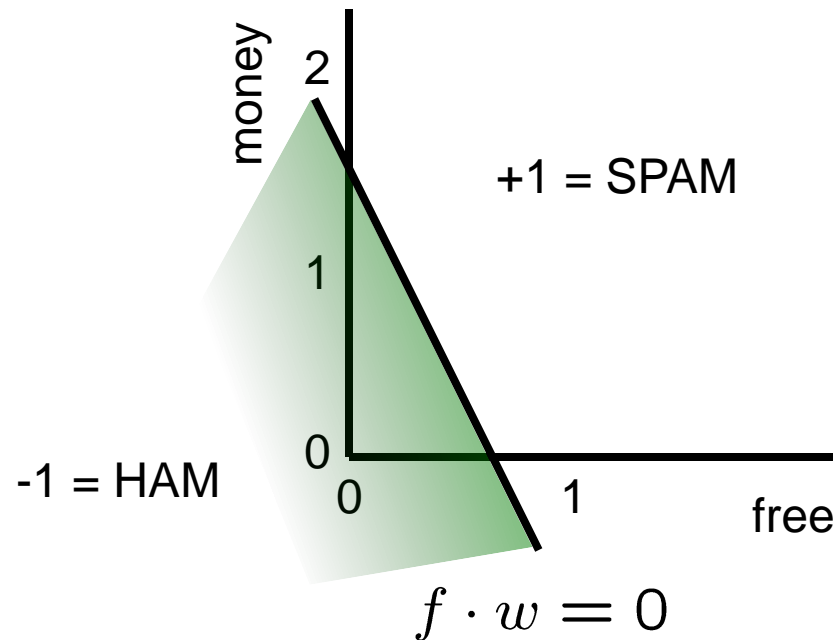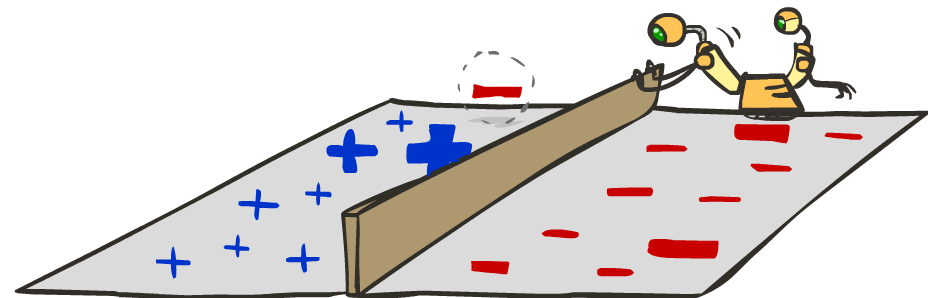  - Other corresponds to Y=-1



$$w$$

```
BIAS  :  -3
free  :   4
money :   2
...
```

+1 = SPAM

-1 = HAM

money

free

$$f \cdot w = 0$$

# 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
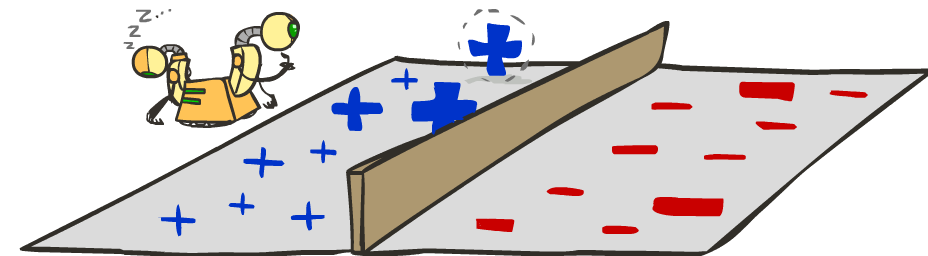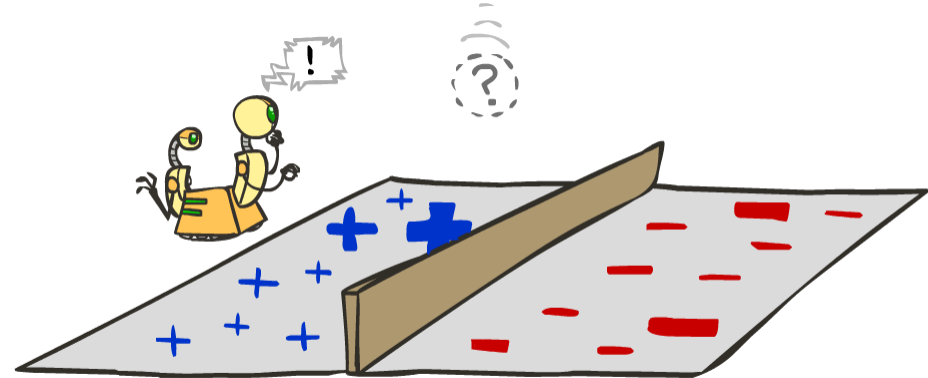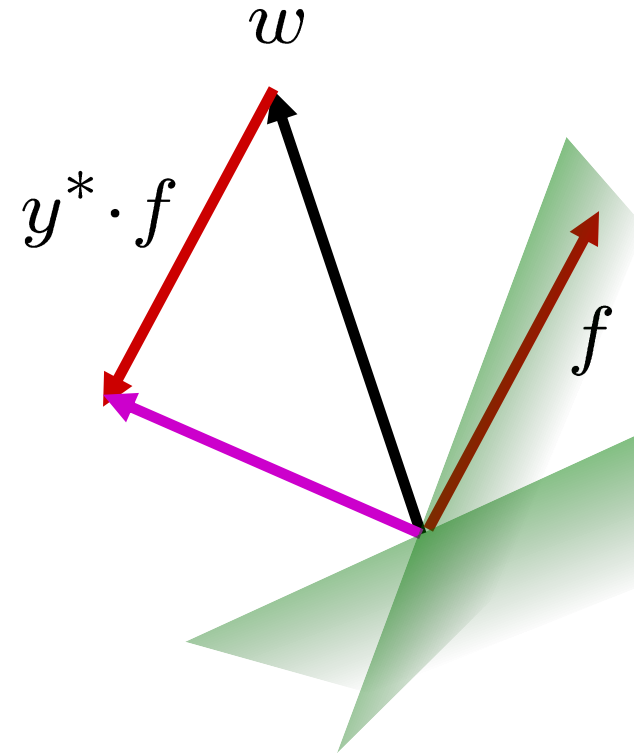
# Learning: Binary Perceptron

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

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

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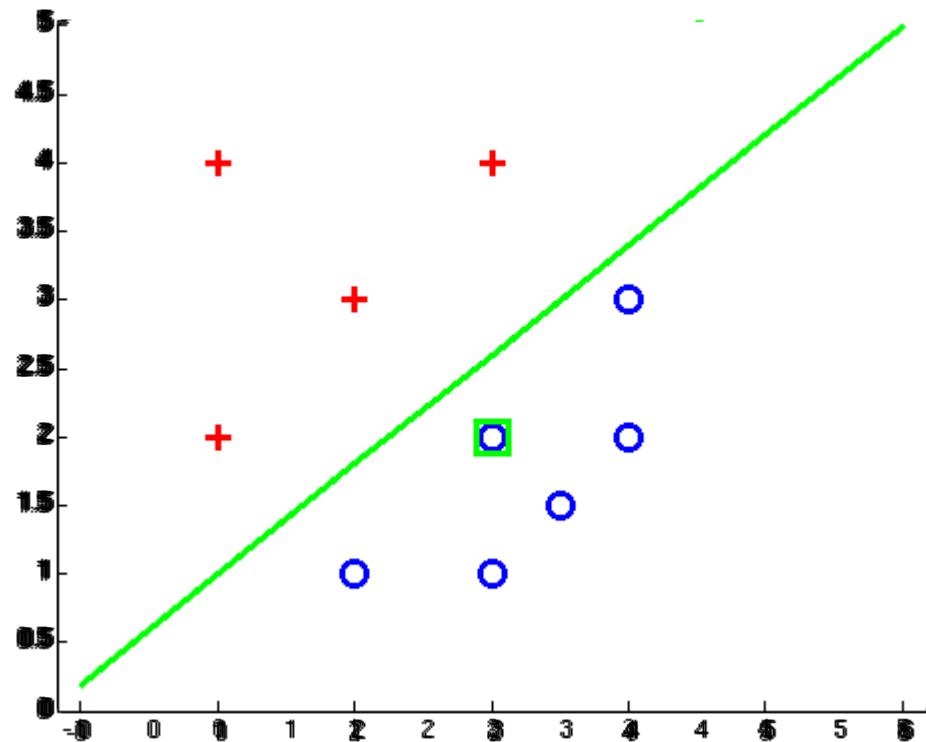  - If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if y* is -1.

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

$w$

$y^* \cdot f$

$f$

Before: *w f*
After: *wf + y\*f f*
*f f >=0*

# Examples: Perceptron

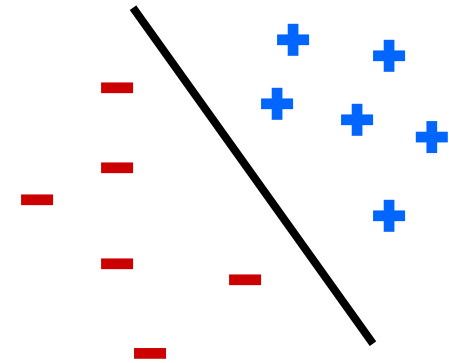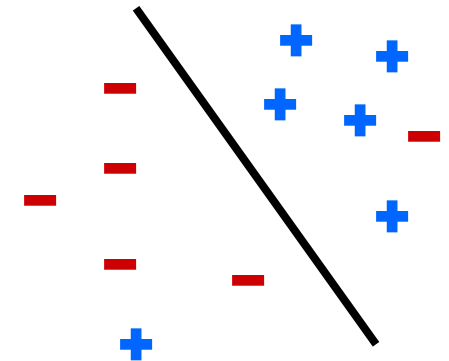- **Separable Case**

# Properties of Perceptrons

- Separability: true if some parameters get the training set perfectly correct

- Convergence: if the training is separable, perceptron will eventually converge (binary case)

- Mistake Bound: the maximum number of mistakes (binary case) related to the *margin* or degree of separability

$$\text{mistakes} < \frac{k}{\delta^2}$$

Separable

Non-Separable

# Problems with the Perceptron

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

- **Mediocre generalization: finds a "barely" separating solution**

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

training

accuracy

test
held-out

iterations

# Improving the Perceptron

# Non-Separable Case: Deterministic Decision


Even the best linear boundary makes at least one mistake

# 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 → want probability going to 1
- If $z = w \cdot f(x)$ very negative → want probability going to 0

- Sigmoid function

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



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

# A 1D Example



$$P(\mathrm{red}|x) = \frac{e^{w_{\mathrm{red}} \cdot x}}{e^{w_{\mathrm{red}} \cdot x} + e^{w_{\mathrm{blue}} \cdot x}}$$

probability increases exponentially as we move away from boundary

normalizer

# Multiclass Logistic Regression

- **Recall Perceptron:**

  - A weight vector for each class: $\quad w_y$

  - Score (activation) of a class y: $\quad w_y \cdot f(x)$

  - Prediction highest score wins $\quad y = \arg\max_y \ w_y \cdot f(x)$



$w_1 \cdot f$ biggest

$w_1$

$w_3$

$w_2$

$w_3 \cdot f$ biggest

$w_2 \cdot f$ biggest

- **How to make the scores into probabilities?**

$$z_1, z_2, z_3 \rightarrow \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$\underbrace{\qquad\qquad}$
original activations

$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad}$
softmax activations

# Best w?

- Maximum likelihood estimation:

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

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**

# Optimization

- i.e., how do we solve:

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

# Hill Climbing

- **Recall from CSPs lecture: 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 $\quad g(w_0 + h) \quad$ and $\quad g(w_0 - h)$

   ■ Then step in best direction

■ Or, evaluate derivative: $\quad \dfrac{\partial g(w_0)}{\partial w} = \lim_{h \to 0} \dfrac{g(w_0 + h) - g(w_0 - h)}{2h}$

   ■ Tells which direction to step into

# 2-D Optimization

# Gradient Ascent

- Perform update in uphill direction for each coordinate

- The steeper the slope (i.e. the higher the derivative) the bigger the step for that coordinate

- E.g., consider: $g(w_1, w_2)$

- Updates:

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

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

- 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
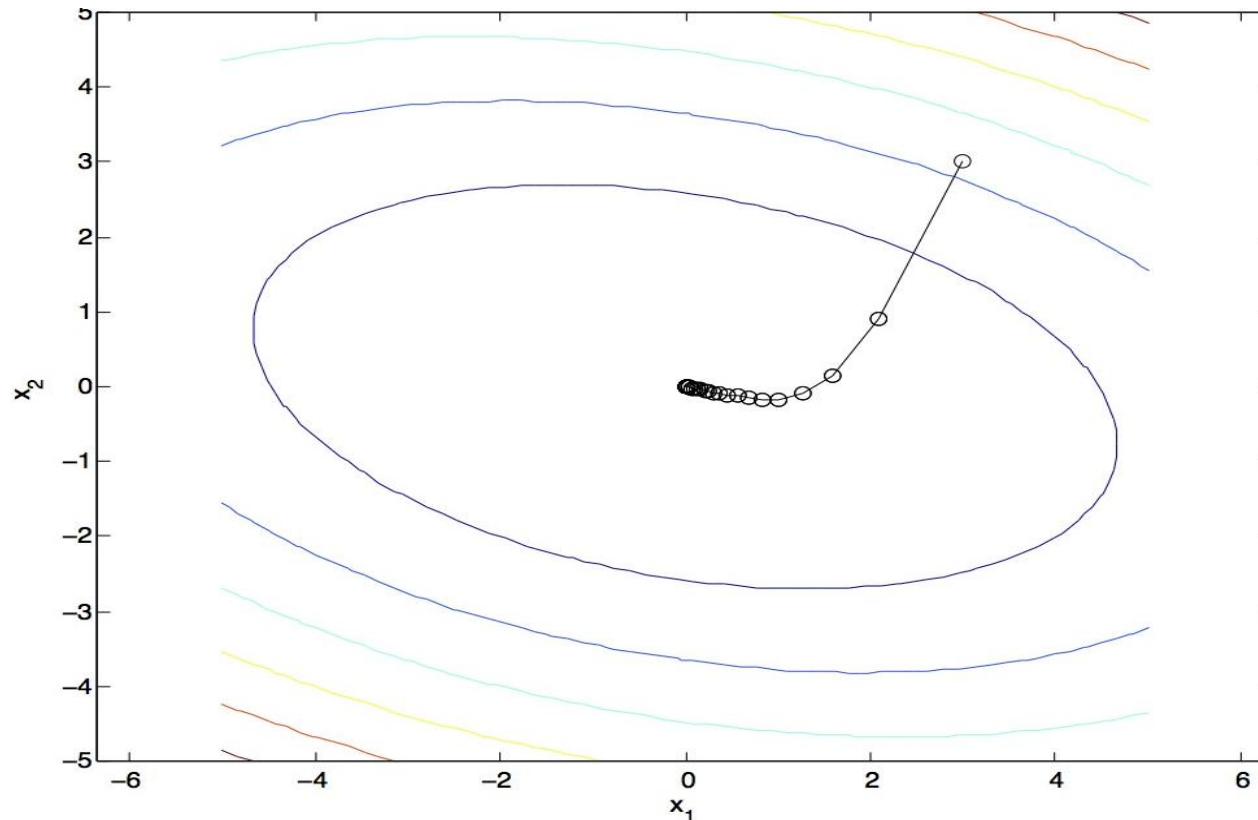
# What is the Steepest Direction?

$$\max_{\Delta:\Delta_1^2+\Delta_2^2\leq\varepsilon} g(w+\Delta)$$

- First-Order Taylor Expansion:

$$g(w+\Delta) \approx g(w) + \frac{\partial g}{\partial w_1}\Delta_1 + \frac{\partial g}{\partial w_2}\Delta_2$$

- Steepest Ascent Direction:

$$\max_{\Delta:\Delta_1^2+\Delta_2^2\leq\varepsilon} g(w) + \frac{\partial g}{\partial w_1}\Delta_1 + \frac{\partial g}{\partial w_2}\Delta_2$$

- Recall:

$$\max_{\Delta:\|\Delta\|\leq\varepsilon} \Delta^\top a \quad \rightarrow \quad \Delta = \varepsilon\frac{a}{\|a\|}$$

- Hence, solution: $\quad \Delta = \varepsilon\dfrac{\nabla g}{\|\nabla g\|}$

**Gradient direction = steepest direction!**

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

# Gradient in n dimensions

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

# Optimization Procedure: Gradient Ascent

- `init` $w$
- `for iter = 1, 2, …`

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

- $\alpha$ : learning rate --- tweaking parameter that needs to be chosen carefully

- How? Try multiple choices

  - Crude rule of thumb: update changes $w$ about $0.1 - 1$ %

# Batch Gradient Ascent on the Log Likelihood Objective

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

$$\underbrace{\phantom{\sum_{i} \log P(y^{(i)}|x^{(i)}; w)}}_{g(w)}$$

- 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

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

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

- `init` $w$
- `for iter = 1, 2, …`
  - `pick random j`
  
  $$w \leftarrow w + \alpha * \nabla \log P(y^{(j)}|x^{(j)}; w)$$

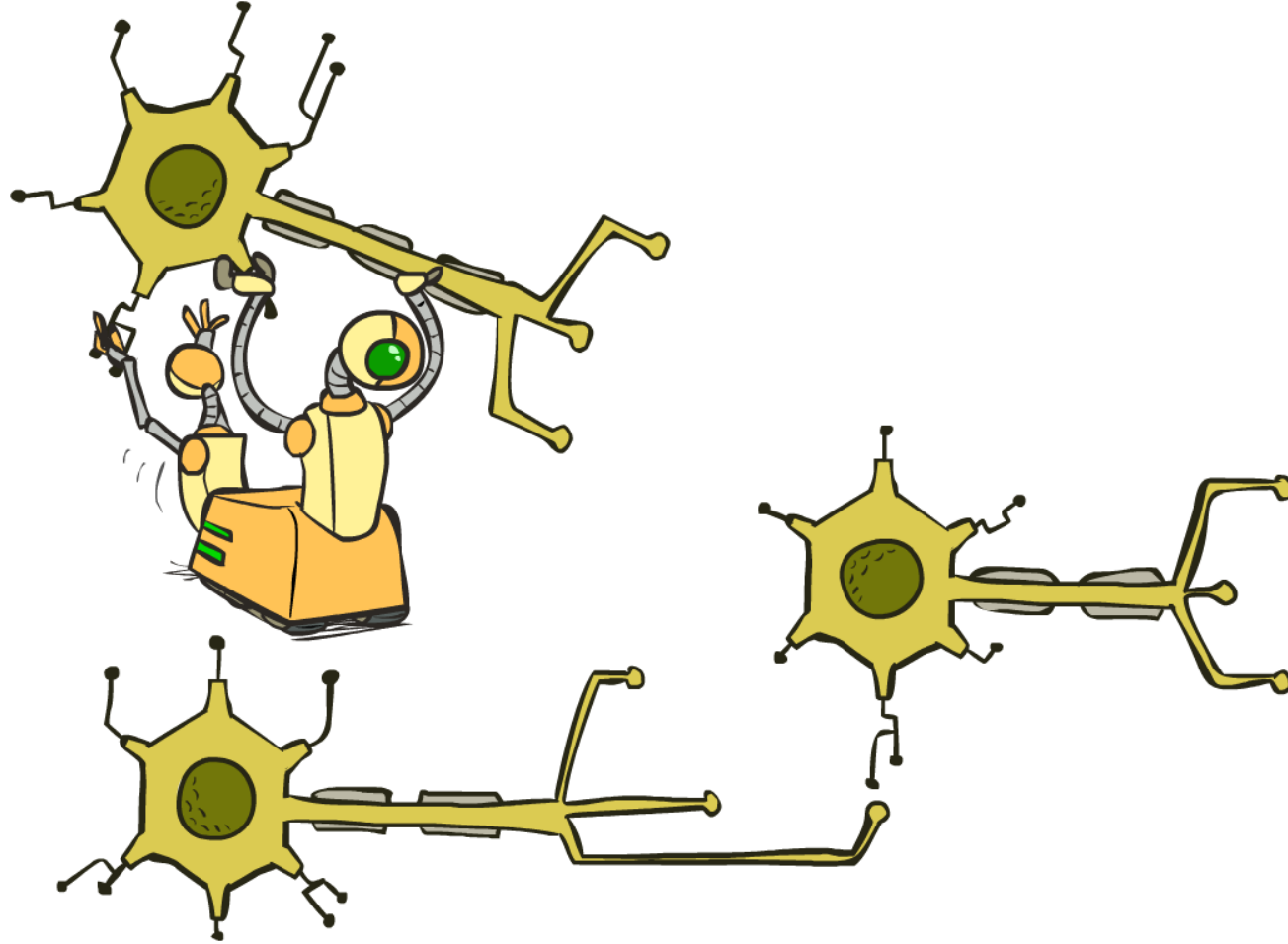# Mini-Batch Gradient Ascent on the Log Likelihood Objective

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

**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

- `init` $w$
- `for iter = 1, 2, …`
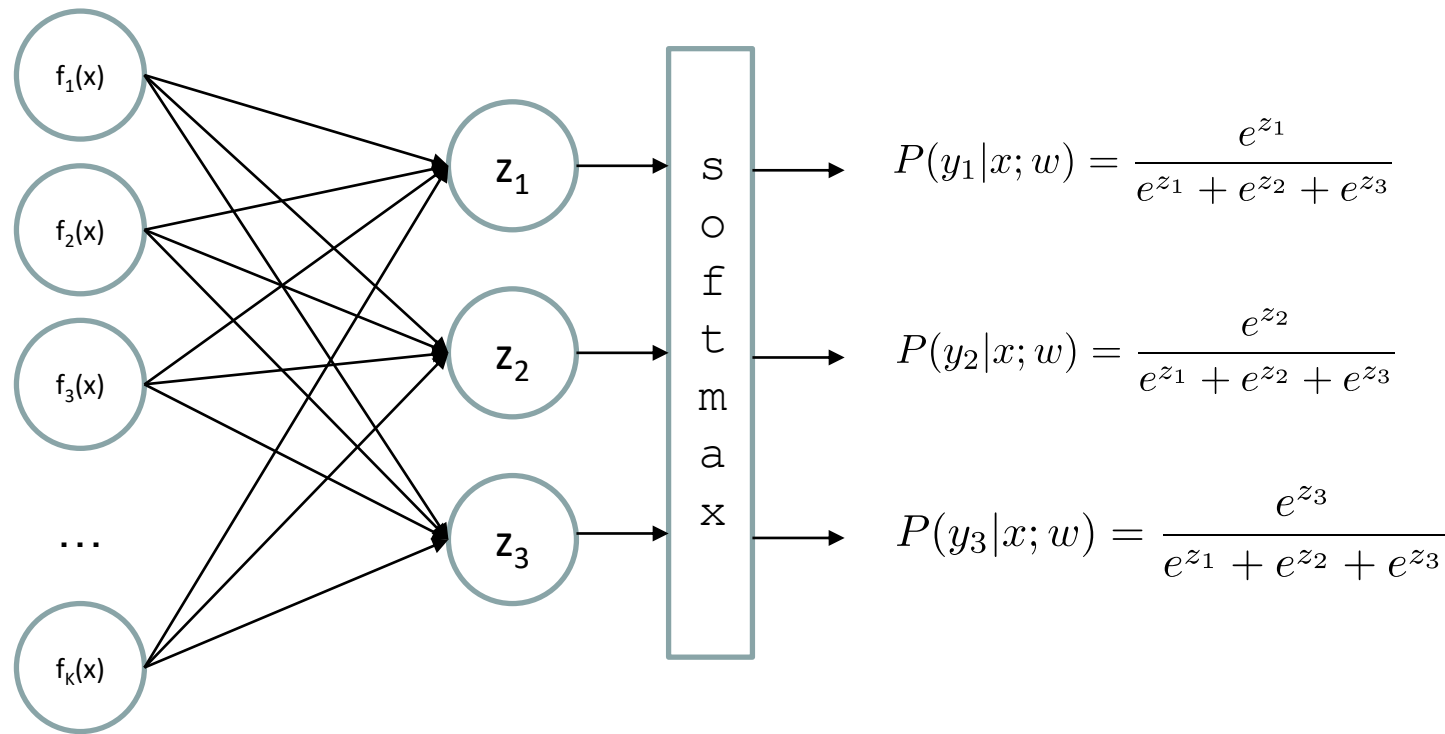  - `pick random subset of training examples J`

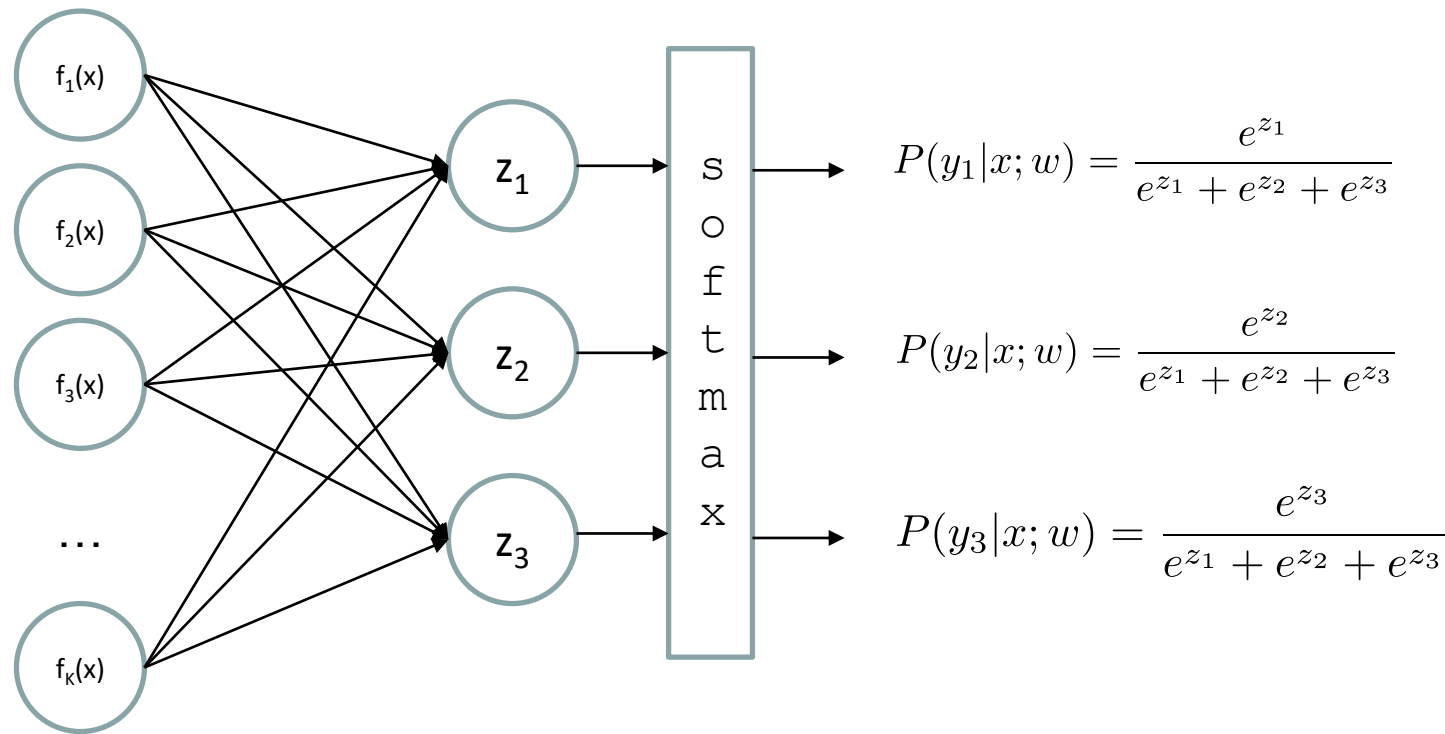  $$w \leftarrow w + \alpha * \sum_{j \in J} \nabla \log P(y^{(j)}|x^{(j)}; w)$$

# Neural Networks

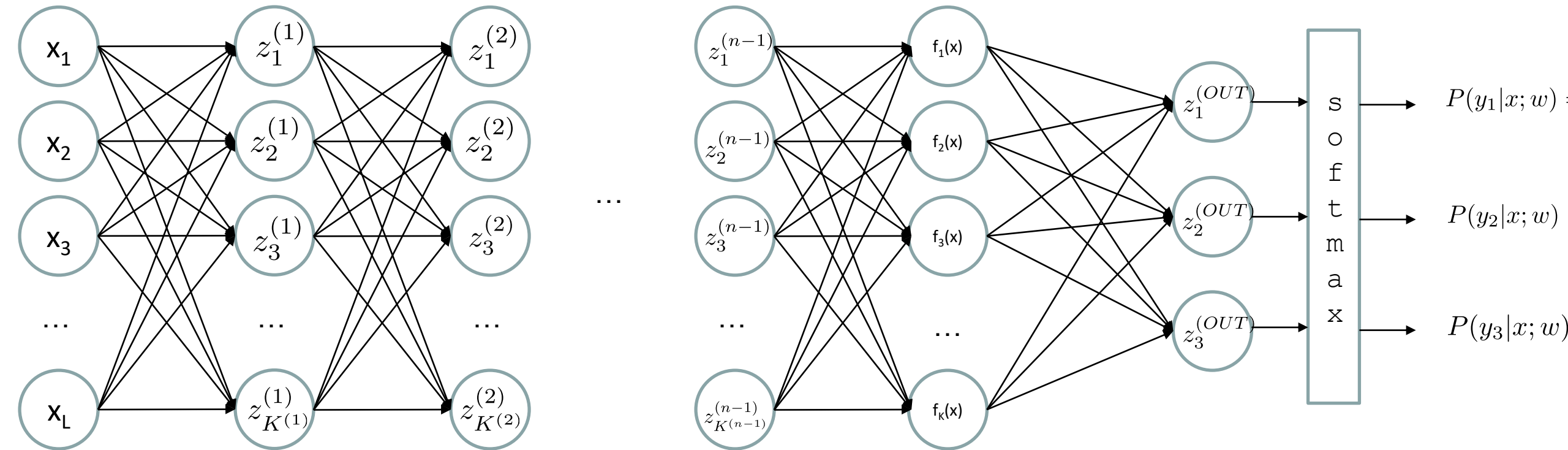# Multi-class Logistic Regression

- = special case of neural network



$$P(y_1|x; w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x; w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3|x; w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# Deep Neural Network = Also learn the features!



$$P(y_1|x;w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x;w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

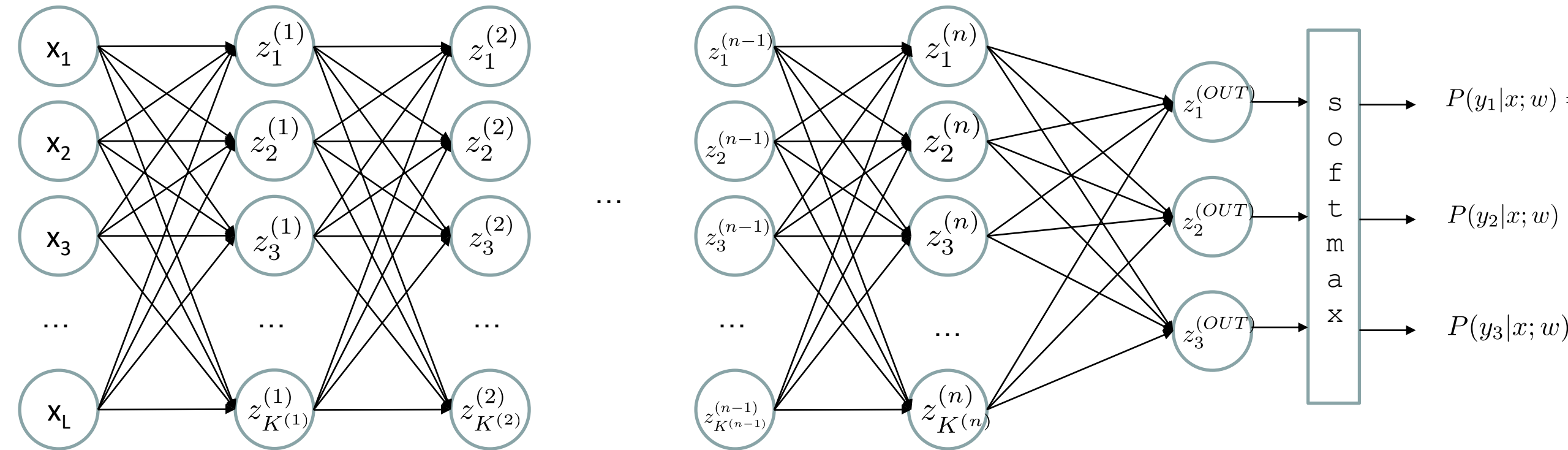$$P(y_3|x;w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# Deep Neural Network = Also learn the features!



$$z_i^{(k)} = g\left(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)}\right)$$

**g = nonlinear activation function**
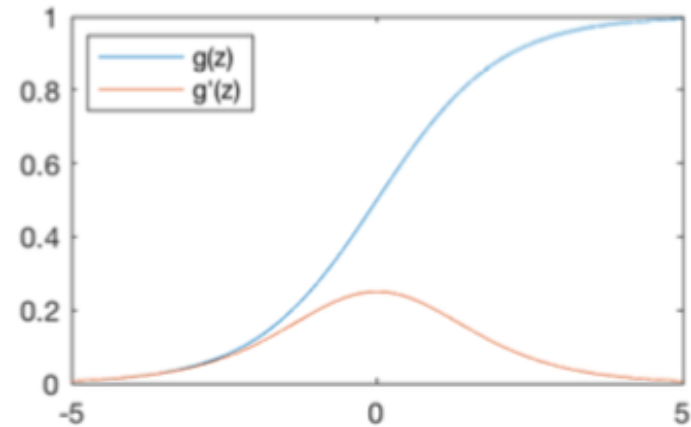
# Deep Neural Network = Also learn the features!



$$z_i^{(k)} = g(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)})$$

**g = nonlinear activation function**
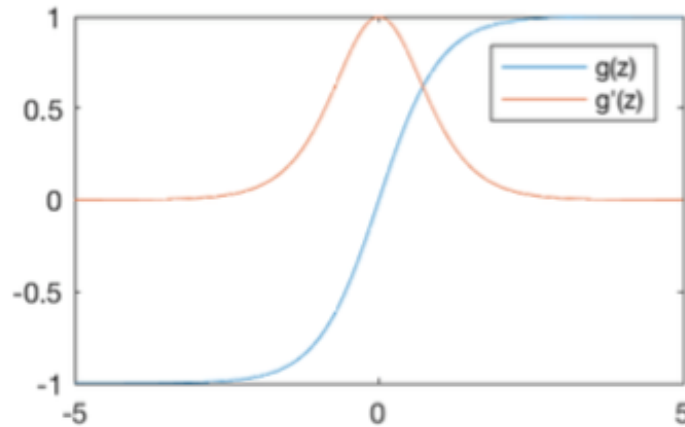
# Common Activation Functions

## Sigmoid Function



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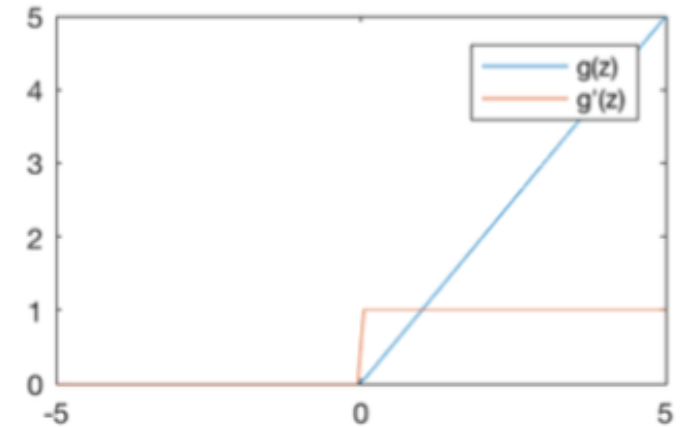
$$g'(z) = g(z)(1 - g(z))$$

## Hyperbolic Tangent



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

$$g'(z) = 1 - g(z)^2$$

## Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

[source: MIT 6.S191 introtodeeplearning.com]

# Deep Neural Network: Also Learn the Features!

- Training the deep neural network is just like logistic regression:

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

just w tends to be a much, much larger vector ☺

→just run gradient ascent

+ stop when log likelihood of hold-out data starts to decrease

# How about computing all the derivatives?

- But neural net f is never one of those?
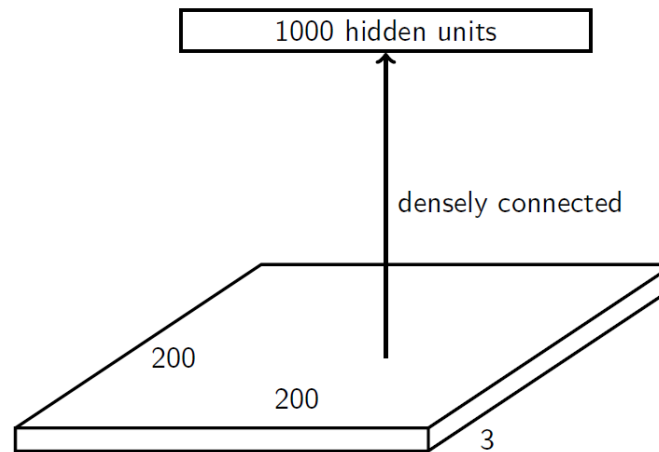  - No problem: CHAIN RULE:

  If $$f(x) = g(h(x))$$

  Then $$f'(x) = g'(h(x))h'(x)$$

  **→ Derivatives can be computed by following well-defined procedures**
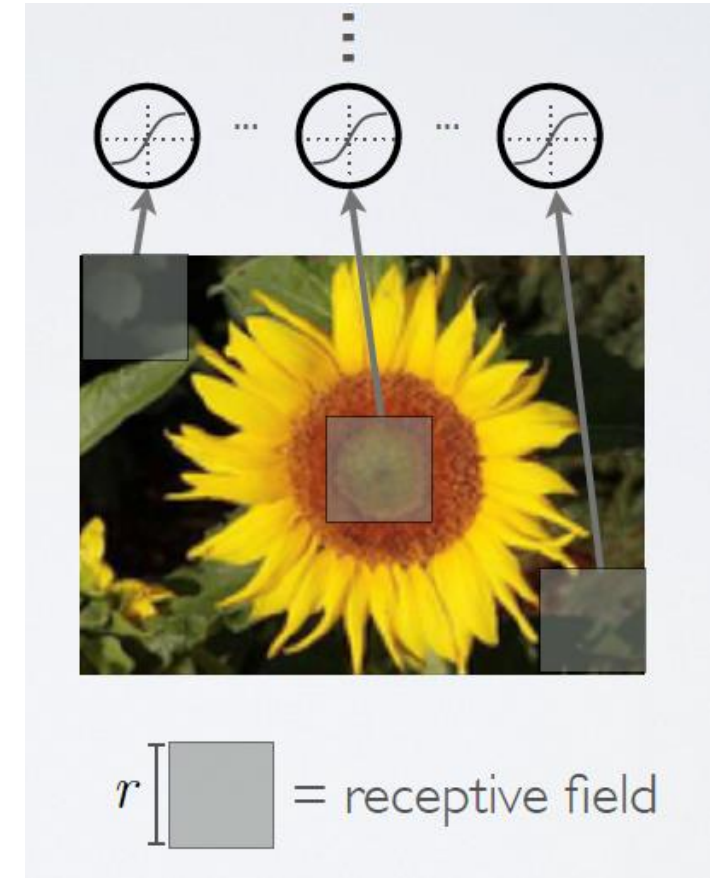
# Motivation

- Visual recognition
  - Suppose we aim to train a network that takes a 200x200 RGB image as input



  - What is the problem with have full connections in the first layer?
    - Too many parameters! 200x200x3x1000 = 120 million
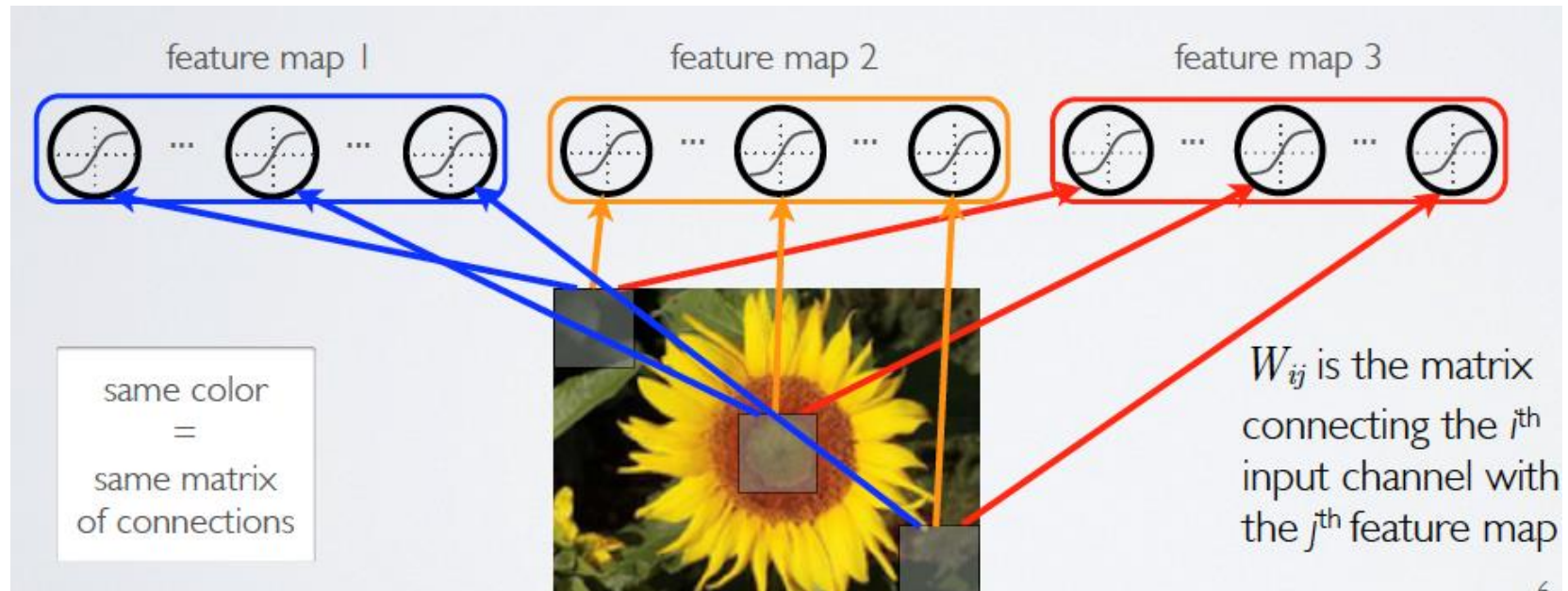    - What happens if the object in the image shifts a little?

# Overview of CNNs

- **First idea: Use a local connectivity of hidden units**
  - Each hidden unit is connected only to a subregion (patch) of the input image
  - Usually it is connected to all channels
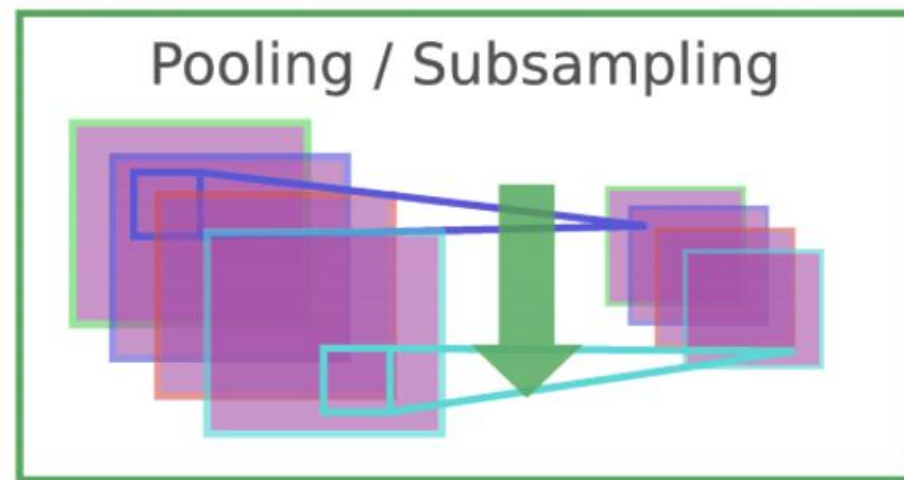  - Each neuron has a local receptive field

# Overview of CNNs

- Second idea: share weights across certain units
  - Units organized into the same "feature map" share weight parameters
  - Hidden units within a feature map cover different positions in the image
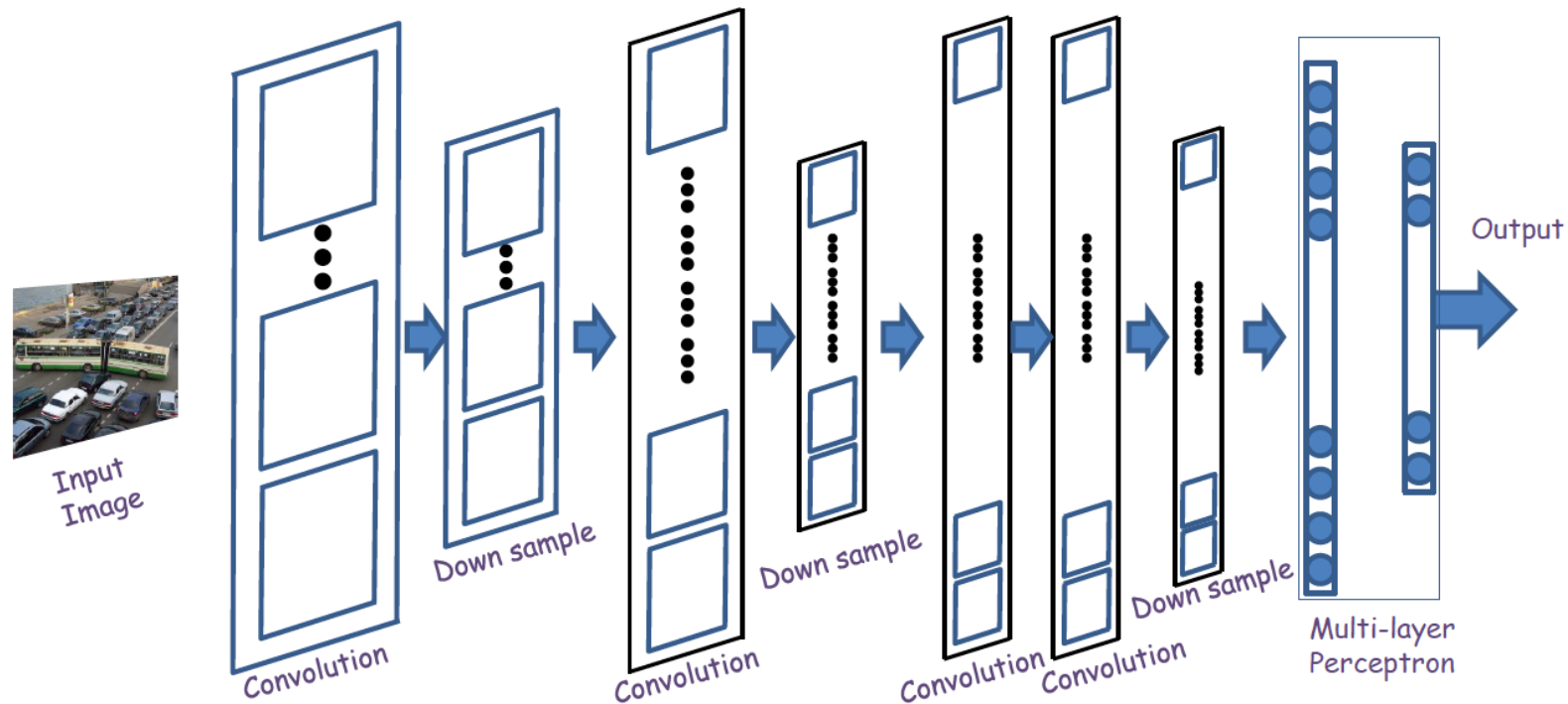
# Overview of CNNs

- Third idea: pool hidden units in the same neighborhood
    - Averaging or Discarding location information in a small region
    - Robust toward small deformations in object shapes by ignoring details.
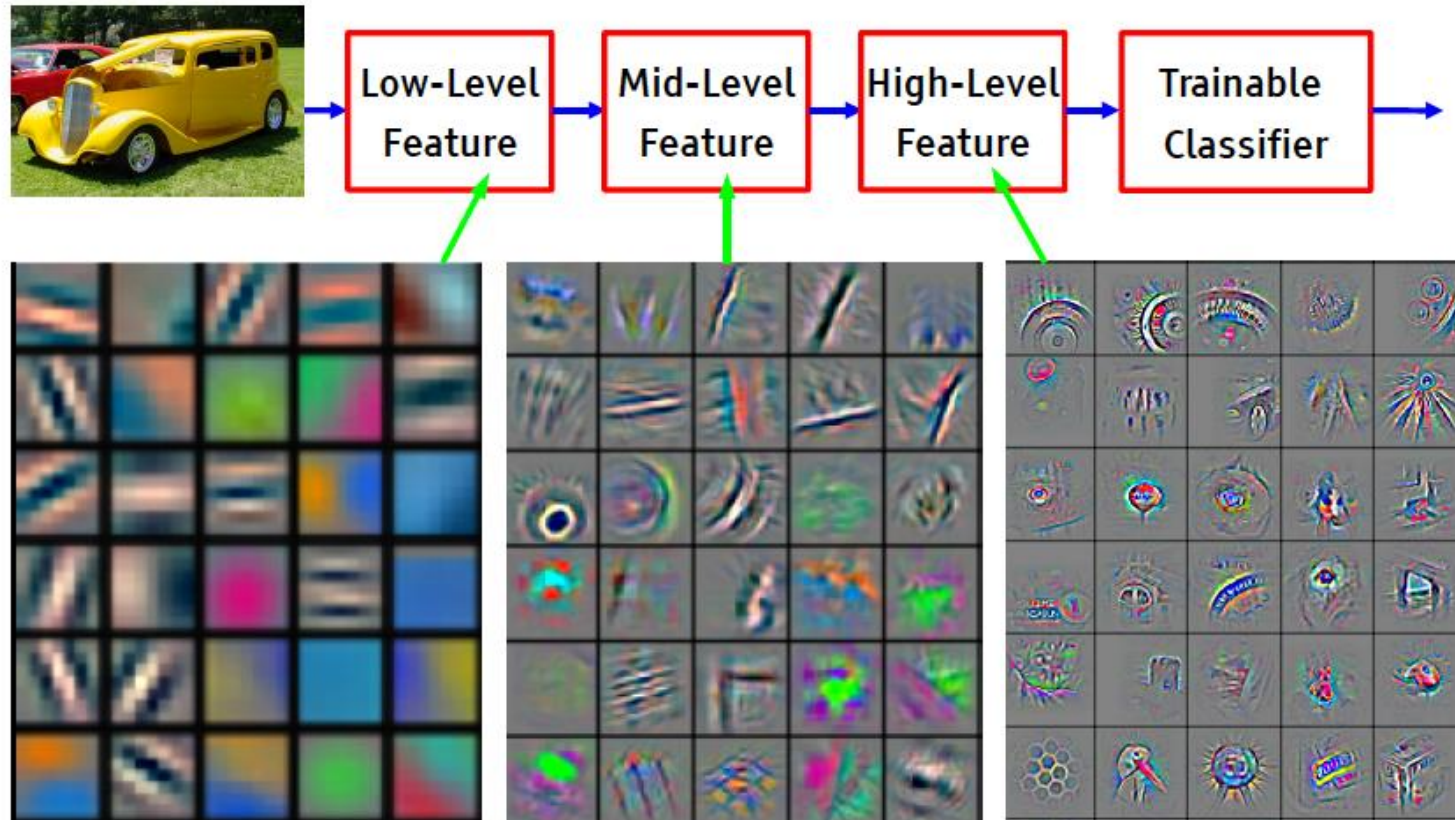
# Overview of CNNs

- Fourth idea: Interleaving feature extraction and pooling operations
    - Extracting abstract, compositional features for representing semantic object classes

# Overview of CNNs

- Artificial visual pathway: from images to semantic concepts (Representation learning)
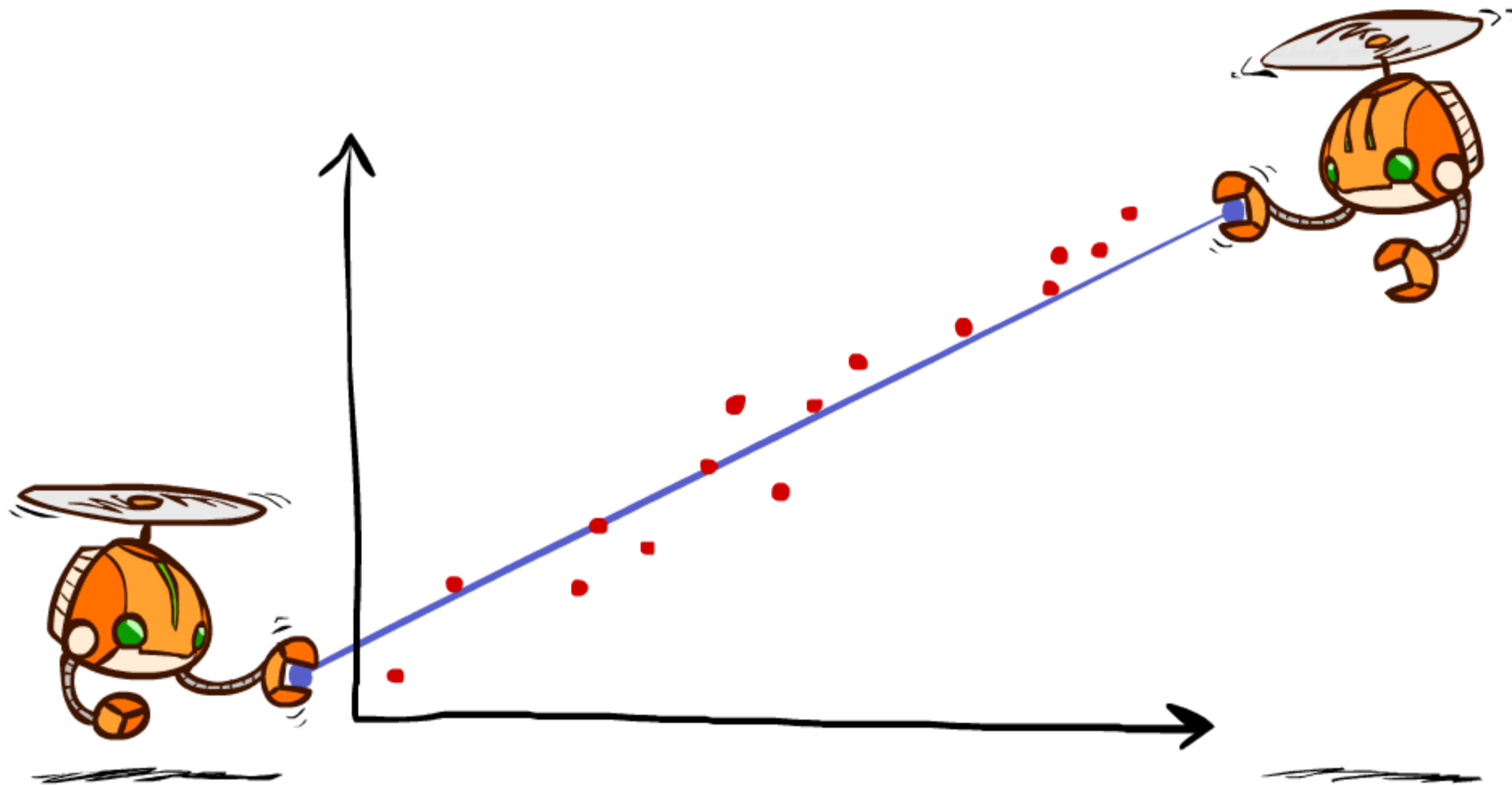


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]
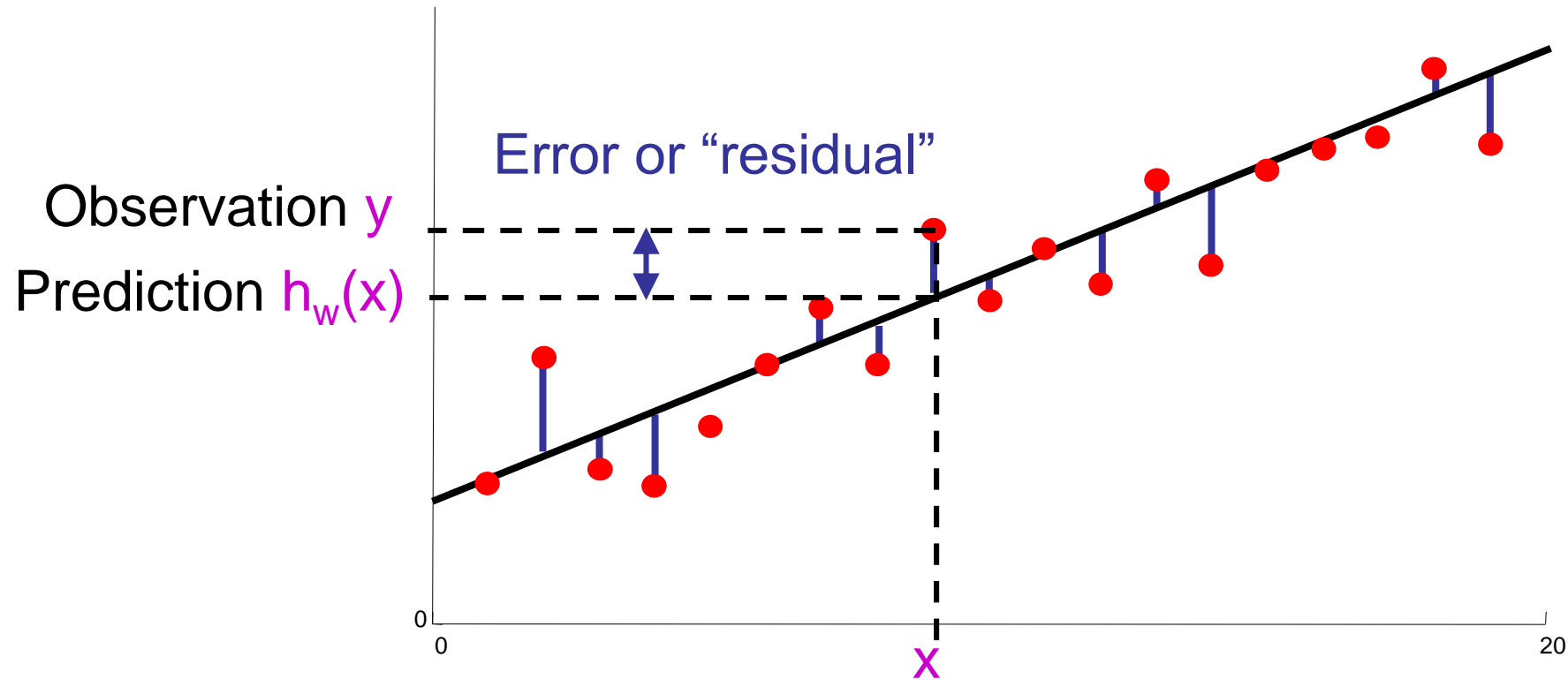
# More classification methods

- Naive Bayes
- Perceptron / Neural networks
- Decision trees / Random forest
- Support Vector Machines
- Nearest neighbors
- Model ensembles: bagging, boosting, etc.
- ......

# Regression

# Linear Regression

Prediction: $h_w(x) = w_0 + w_1 x$



Error or "residual"

Observation $y$

Prediction $h_w(x)$

Error on one instance: $|y - h_w(x)|$

# Least squares: Minimizing squared error

- L2 loss function: sum of squared errors over all examples

$$L(\boldsymbol{w}) = \sum_i \left(y_i - h_w(\boldsymbol{x}_i)\right)^2 = \sum_i (y_i - \boldsymbol{w}^T \boldsymbol{x}_i)^2$$

- We want the weights w* that minimize loss
- Analytical solution: at w* the derivative of loss w.r.t. each weight is zero
  - **X** is the data matrix (all the data, one example per row); **y** is the vector of labels
  - **w*** = (**X**$^T$**X**)$^{-1}$**X**$^T$**y**

# Regularized Regression

- Overfitting is also possible in regression
  - Extreme case: $n$ features, $n$ training examples
- Regularization can be used to alleviate overfitting

- LASSO (Least Absolute Shrinkage and Selection Operator)

$$L(\boldsymbol{w}) = \sum_i (y_i - \boldsymbol{w}^T \boldsymbol{x}_i)^2 + \lambda \sum_k |w_k|$$

- Ridge Regression

$$L(\boldsymbol{w}) = \sum_i (y_i - \boldsymbol{w}^T \boldsymbol{x}_i)^2 + \lambda \sum_k w_k^2$$