# 优化方法与神经网络
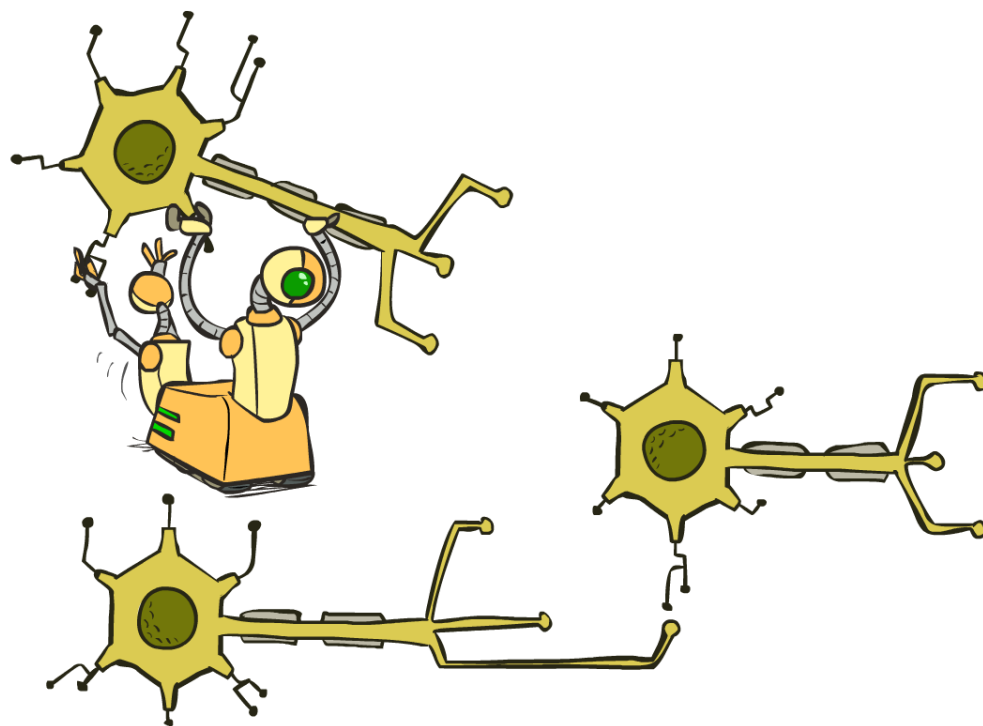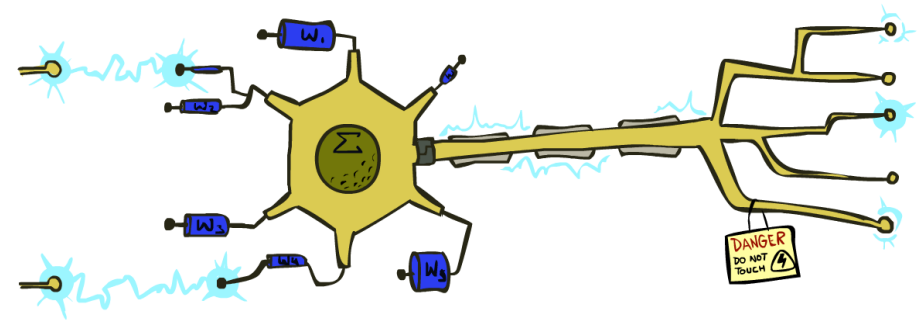
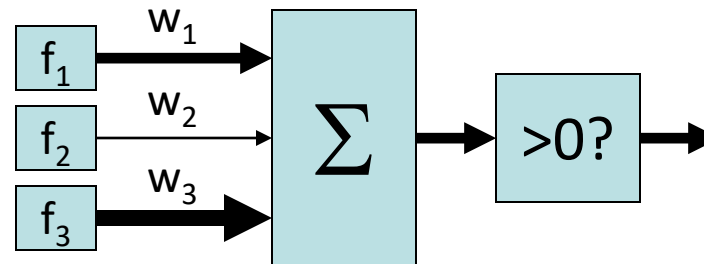# 回顾: 线性判别分类器

- 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:
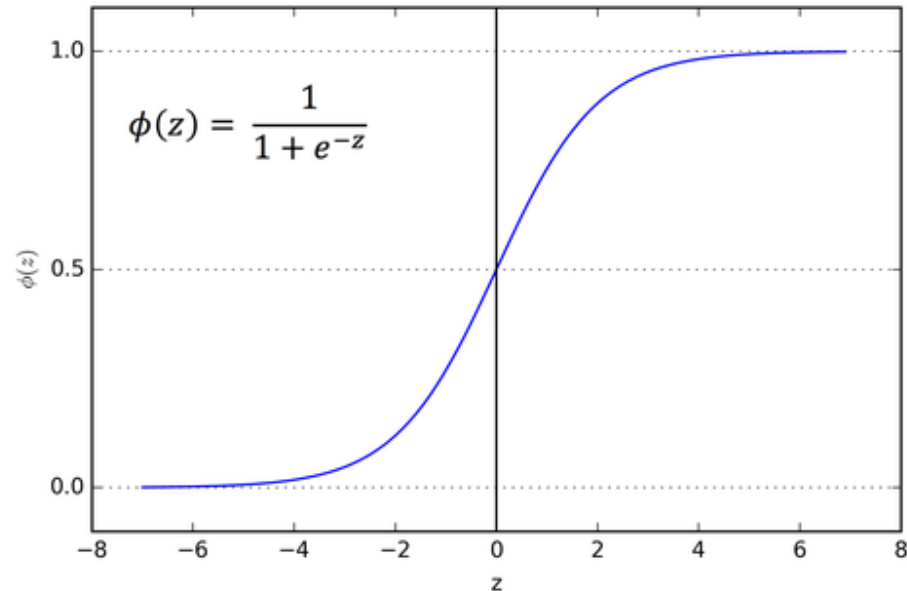  - Positive, output +1
  - Negative, output -1

# 如何获得概率化的判别决策?

- Activation: $z = w \cdot f(x)$
- 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

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



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

# 求解最优的 w?

- Maximum likelihood estimation:

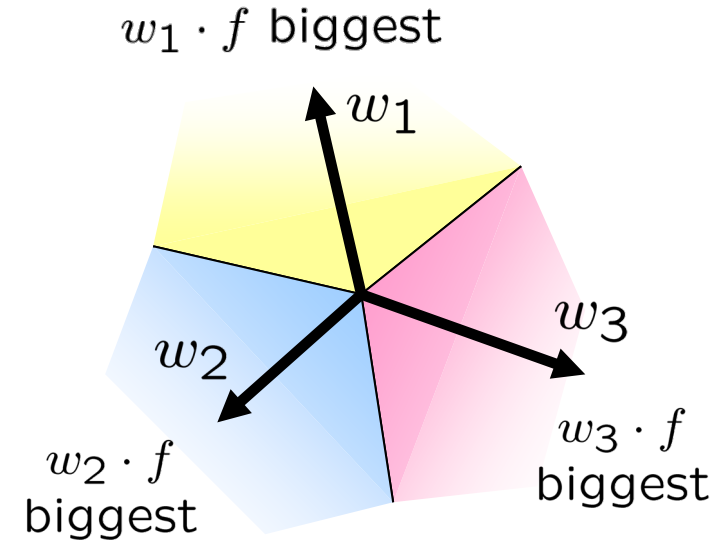$$\max_w \quad ll(w) = \max_w \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

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**

# 多分类罗吉斯特回归

- **Multi-class linear classification**

  - A weight vector for each class: $w_y$

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

  - Prediction w/highest score wins: $y = \underset{y}{\arg\max} \; w_y \cdot f(x)$

$w_1 \cdot f$ biggest

$w_1$

$w_3$

$w_2$

$w_2 \cdot f$
biggest

$w_3 \cdot f$
biggest

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

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

original activations

softmax activations

# 求解最优的 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**

# 优化问题

- Optimization

  - i.e., how do we solve:

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

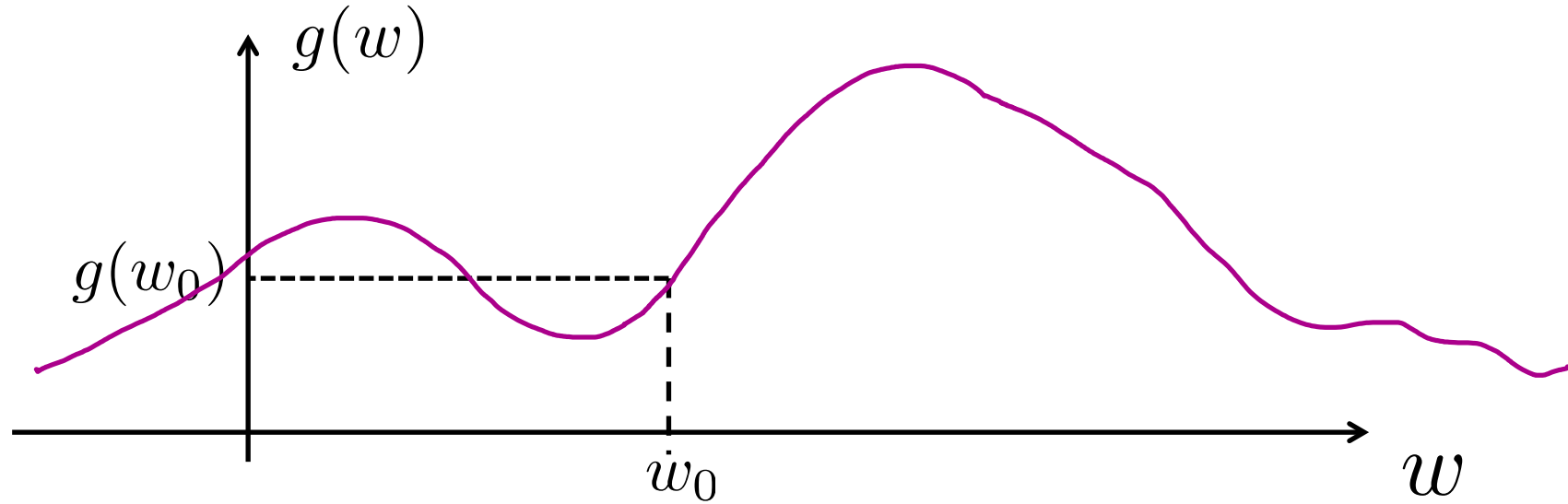# Hill Climbing 爬山算法

- **在约束满足问题里面介绍过: simple, general idea**
  - Start wherever
  - Repeat: move to the best neighboring state
  - If no neighbors better than current, quit

- **这里的挑战，求解多分类下罗吉斯特回归优化问题?**
  - Optimization over a continuous space 连续空间
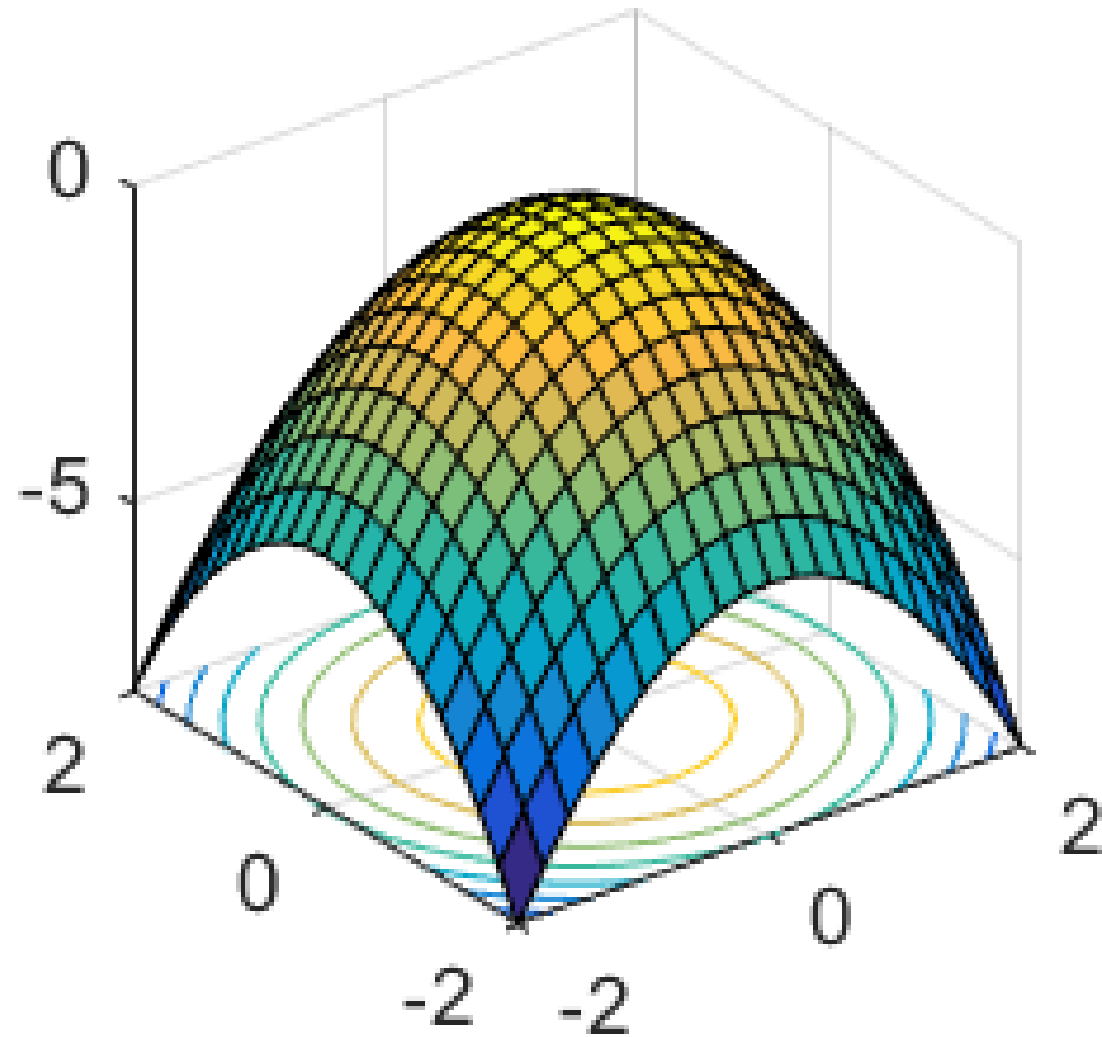    - Infinitely many neighbors!
    - How to do this efficiently?

# 一维优化



- Could evaluate $g(w_0 + h)$ and $g(w_0 - h)$

  - Then step in best direction

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

  - Tells which direction to step into

# 2-D Optimization

# Gradient Ascent 梯度升高法

- 把每一维度的权值推向上山的方向
- 梯度越陡 (i.e. 导数越大) 更新的步长就越大
- 例如:

$$g(w_1, w_2)$$

- Updates:

- Updates in vector notation:

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

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

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

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

# 梯度升高法

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

# 求解最陡的方向?

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

- 一阶泰勒展开: $$g(w+\Delta) \approx g(w) + \frac{\partial g}{\partial w_1}\Delta_1 + \frac{\partial g}{\partial w_2}\Delta_2$$

- 最陡峭的爬升方向: $$\max_{\Delta:\Delta_1^2+\Delta_2^2\leq\varepsilon} g(w) + \frac{\partial g}{\partial w_1}\Delta_1 + \frac{\partial g}{\partial w_2}\Delta_2$$

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

- Hence, solution: $\Delta = \varepsilon\frac{\nabla g}{\|\nabla g\|}$ **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 梯度

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

# 优化过程: 梯度上升法

- `init` $w$
- `for iter = 1, 2, …`

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

- $\alpha$: 学习率 --- tweaking parameter that needs to be chosen carefully
- How? Try multiple choices
  - 经验做法: update changes $w$ about $0.1 - 1$ %

# Batch Gradient Ascent on the Log Likelihood Objective

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

- `init` $w$
- `for iter = 1, 2, …`

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

# 在梯度上升法中每个权值向量的更新?

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

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

$$\nabla w_{y^{(i)}} f(x^{(i)}) - \nabla \log \sum_y e^{w_y f(x^{(i)})}$$

adds f to the correct
class weights

$$\frac{1}{\sum_y e^{w_y f(x^{(i)})}} \sum_y \left( e^{w_y f(x^{(i)})} [0^T f(x^{(i)})^T 0^T]^T \right)$$

for y' weights: $\frac{1}{\sum_y e^{w_y f(x^{(i)})}} e^{w_{y'} f(x^{(i)})} f(x^{(i)})$

$$P(y'|x^{(i)}; w) f(x^{(i)})$$

subtracts f from y' weights in proportion to
the probability current weights give to y'

# 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

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

# How about computing all the derivatives?

- We'll talk about that once we covered neural networks, which are a generalization of logistic regression

# Neural Networks 神经网络

# 多分类罗吉斯特回归

- = 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深度神经网络 = 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**

# 常用的激活函数



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)

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

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

# 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

# 神经网络的属性

- **Theorem (Universal Function Approximators). A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.**

- **Practical considerations**
  - Can be seen as learning the features

  - Large number of neurons
    - Danger for overfitting
    - (hence early stopping!)

# Universal Function Approximation Theorem*

**Hornik theorem 1:** Whenever the activation function is *bounded and nonconstant*, then, for any finite measure $\mu$, standard multilayer feedforward networks can approximate any function in $L^p(\mu)$ (the space of all functions on $R^k$ such that $\int_{R^k} |f(x)|^p d\mu(x) < \infty$) arbitrarily well, provided that sufficiently many hidden units are available.

**Hornik theorem 2:** Whenever the activation function is *continuous, bounded and nonconstant*, then, for arbitrary compact subsets $X \subseteq R^k$, standard multilayer feedforward networks can approximate any continuous function on $X$ arbitrarily well with respect to uniform distance, provided that sufficiently many hidden units are available.

- <u>In words:</u> Given any continuous function f(x), if a 2-layer neural network has enough hidden units, then there is a choice of weights that allow it to closely approximate f(x).

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"
Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"
Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

# Universal Function Approximation Theorem*

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"
Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"
Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

# 神经网络演示网址

- Demo-site:
  - http://playground.tensorflow.org/

# How about computing all the derivatives（求导函数）？

- Derivatives tables:

$$\frac{d}{dx}(a) = 0$$

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(au) = a\frac{du}{dx}$$

$$\frac{d}{dx}(u+v-w) = \frac{du}{dx} + \frac{dv}{dx} - \frac{dw}{dx}$$

$$\frac{d}{dx}(uv) = u\frac{dv}{dx} + v\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{u}{v}\right) = \frac{1}{v}\frac{du}{dx} - \frac{u}{v^2}\frac{dv}{dx}$$

$$\frac{d}{dx}(u^n) = nu^{n-1}\frac{du}{dx}$$

$$\frac{d}{dx}(\sqrt{u}) = \frac{1}{2\sqrt{u}}\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u}\right) = -\frac{1}{u^2}\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u^n}\right) = -\frac{n}{u^{n+1}}\frac{du}{dx}$$

$$\frac{d}{dx}[f(u)] = \frac{d}{du}[f(u)]\frac{du}{dx}$$

$$\frac{d}{dx}[\ln u] = \frac{d}{dx}[\log_e u] = \frac{1}{u}\frac{du}{dx}$$

$$\frac{d}{dx}[\log_a u] = \log_a e\frac{1}{u}\frac{du}{dx}$$

$$\frac{d}{dx}e^u = e^u\frac{du}{dx}$$

$$\frac{d}{dx}a^u = a^u\ln a\frac{du}{dx}$$

$$\frac{d}{dx}(u^v) = vu^{v-1}\frac{du}{dx} + \ln u \ u^v\frac{dv}{dx}$$

$$\frac{d}{dx}\sin u = \cos u\frac{du}{dx}$$

$$\frac{d}{dx}\cos u = -\sin u\frac{du}{dx}$$

$$\frac{d}{dx}\tan u = \sec^2 u\frac{du}{dx}$$

$$\frac{d}{dx}\cot u = -\csc^2 u\frac{du}{dx}$$

$$\frac{d}{dx}\sec u = \sec u\tan u\frac{du}{dx}$$

$$\frac{d}{dx}\csc u = -\csc u\cot u\frac{du}{dx}$$

# How about computing all the derivatives?

- But neural net f is never one of those?
  - No problem: CHAIN RULE（求导链式法则）：

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

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

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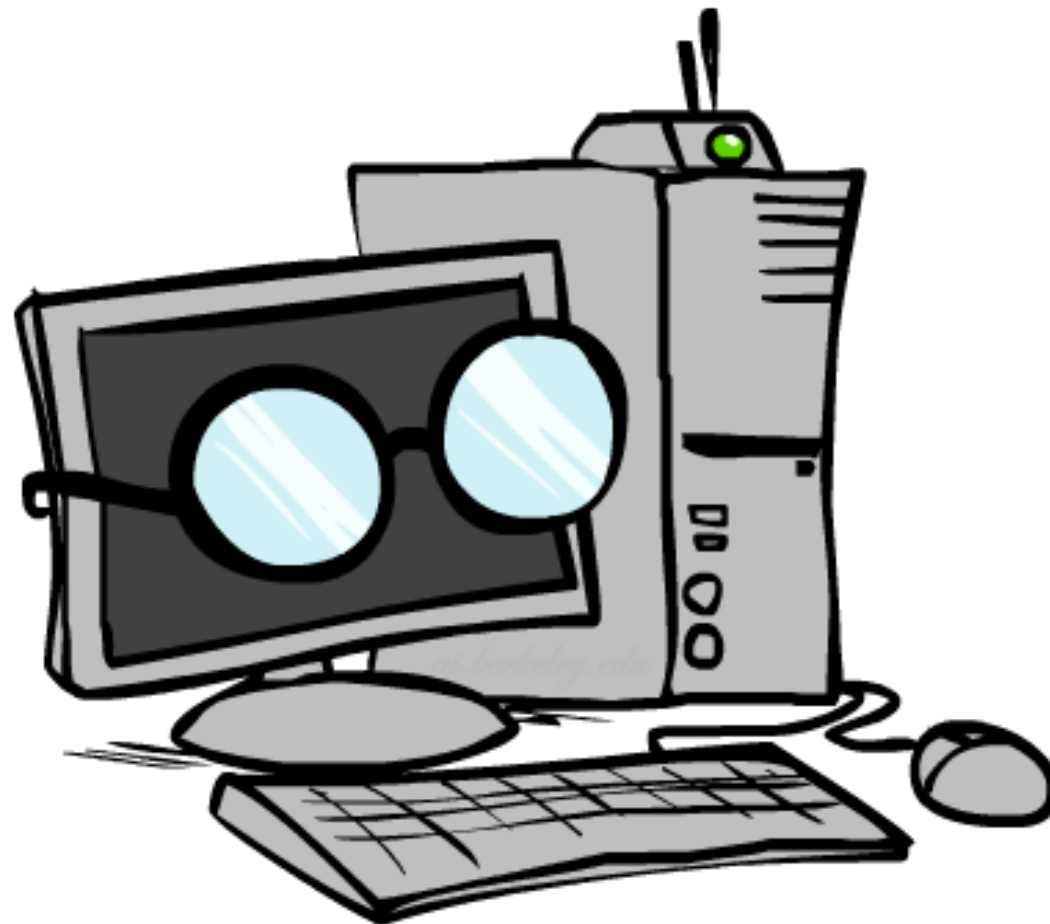
# Automatic Differentiation

- **Automatic differentiation software**
  - e.g. Theano, TensorFlow, PyTorch, Chainer
  - Only need to program the function g(x,y,w)
  - Can automatically compute all derivatives w.r.t. all entries in w
  - This is typically done by caching info during forward computation pass of f, and then doing a backward pass = "backpropagation"
  - Autodiff / Backpropagation can often be done at computational cost comparable to the forward pass
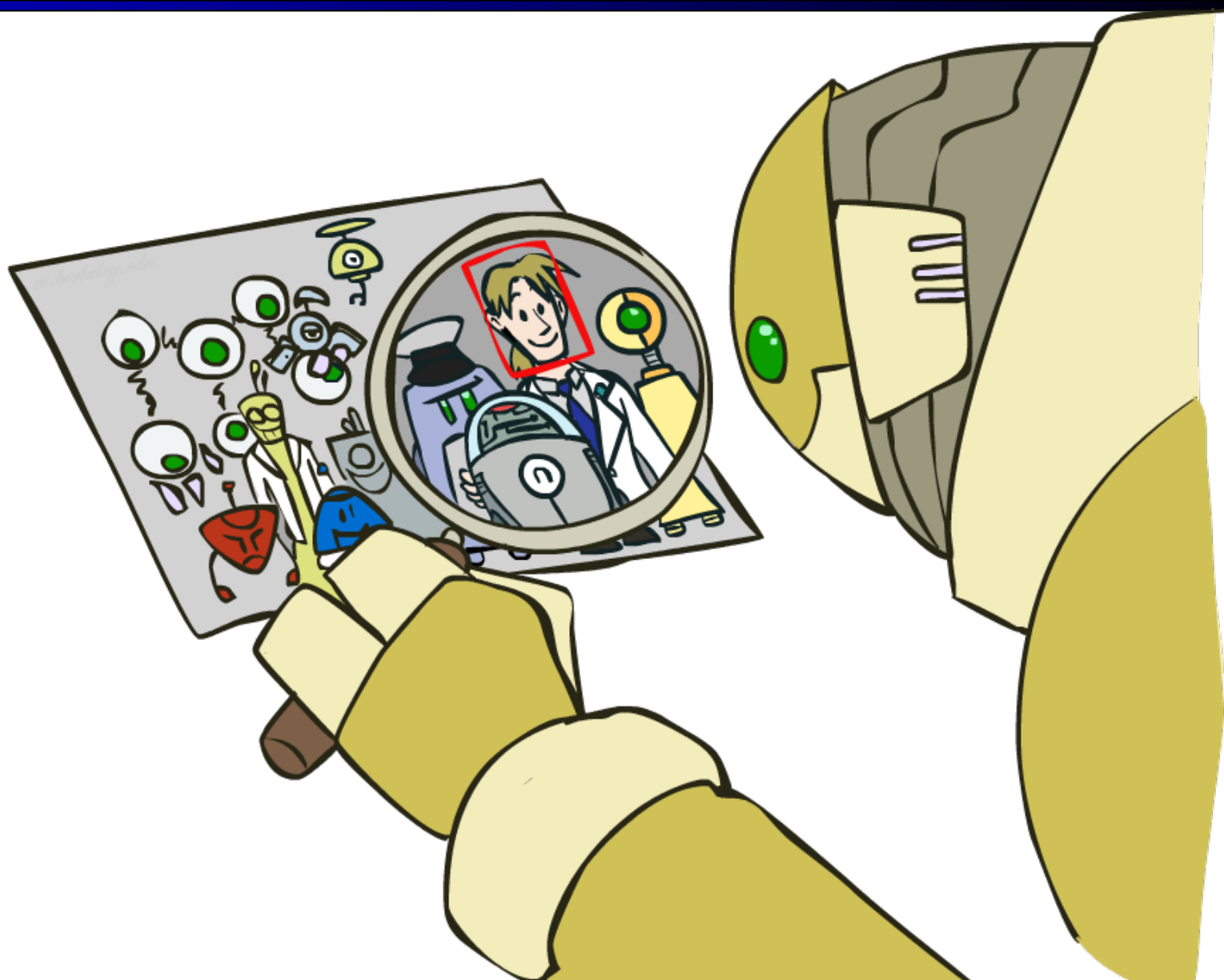- **Need to know this exists**
- **How this is done?**

# 小结

- Optimize probability of label given input $$\max_w \; ll(w) = \max_w \; \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

- Continuous optimization
  - Gradient ascent:
    - Compute steepest uphill direction = gradient (= just vector of partial derivatives)
    - Take step in the gradient direction
    - Repeat (until held-out data accuracy starts to drop = "early stopping")

- Deep neural nets
  - Last layer = still logistic regression
  - Now also many more layers before this last layer
    - = computing the features
    - → the features are learned rather than hand-designed
  - Universal function approximation theorem
    - `If` neural net is large enough
    - `Then` neural net can represent any continuous mapping from input to output with arbitrary accuracy
    - But remember: need to avoid overfitting / memorizing the training data → early stopping!
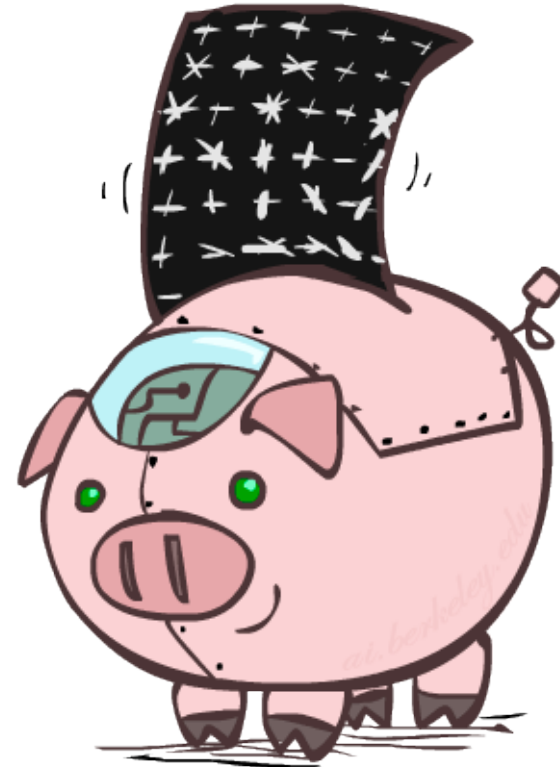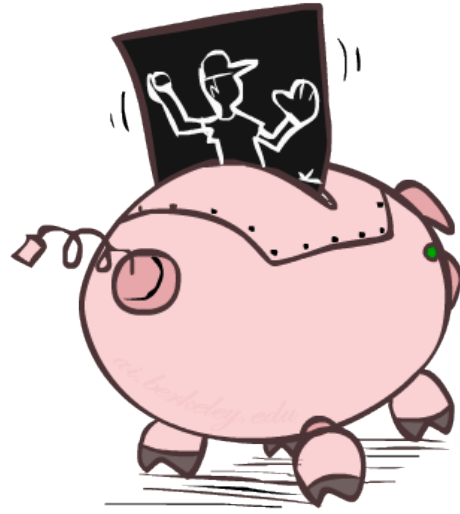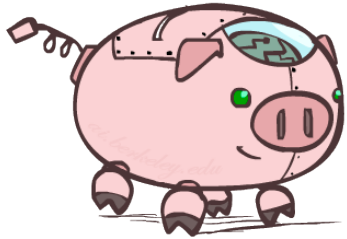  - Automatic differentiation gives the derivatives efficiently
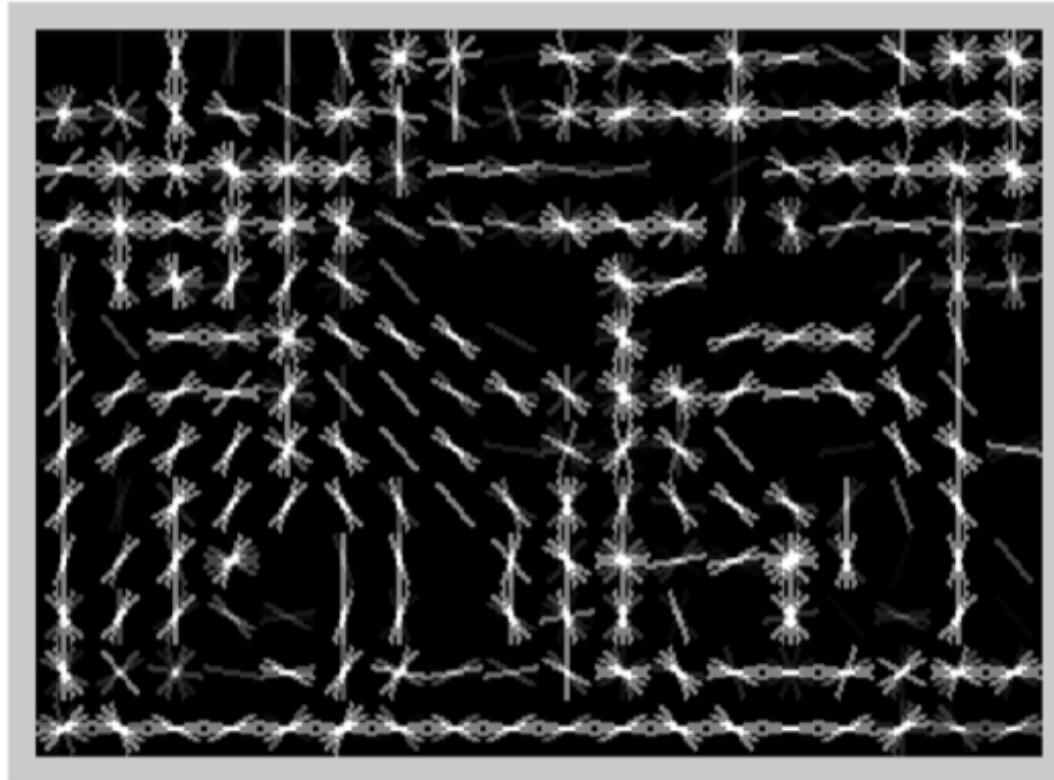
# Computer Vision 计算机视觉

# Manual Feature Design

# Features and Generalization



[HoG: Dalal and Triggs, 2005]

# Features and Generalization



Image

HoG

# Performance



*graph credit Matt Zeiler, Clarifai*

# Performance



## ImageNet Error Rate 2010-2014

*graph credit Matt Zeiler, Clarifai*

# Performance



ImageNet Error Rate 2010-2014

*graph credit Matt Zeiler, Clarifai*

# Performance



*graph credit Matt Zeiler, Clarifai*

# Performance



ImageNet Error Rate 2010-2014

*graph credit Matt Zeiler, Clarifai*
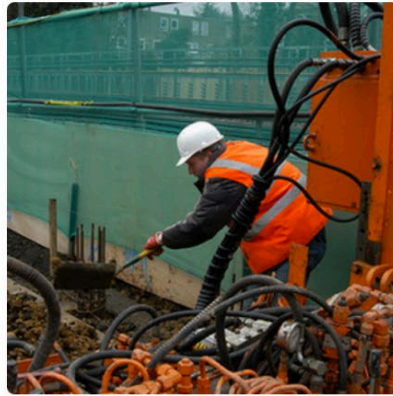
# MS COCO Image Captioning Challenge



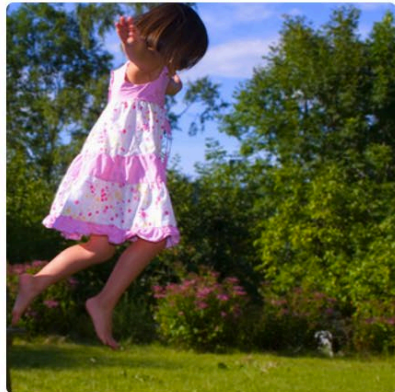"man in black shirt is playing guitar."

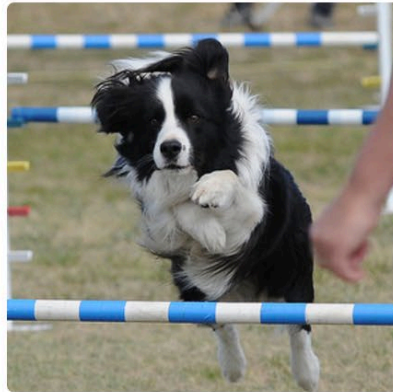"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

"girl in pink dress is jumping in air."

"black and white dog jumps over bar."

"young girl in pink shirt is swinging on swing."

"man in blue wetsuit is surfing on wave."

Karpathy & Fei-Fei, 2015; Donahue et al., 2015; Xu et al, 2015; many more

# Visual QA Challenge

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh

# Speech Recognition



*graph credit Matt Zeiler, Clarifai*

# Machine Translation

(a) Husky classified as wolf    (b) Explanation

Figure 11: Raw data and explanation of a bad model's prediction in the "Husky vs Wolf" task.

|  | Before | After |
|---|---|---|
| Trusted the bad model | 10 out of 27 | 3 out of 27 |
| Snow as a potential feature | 12 out of 27 | 25 out of 27 |

Table 2: "Husky vs Wolf" experiment results.

[Ribeiro et al.]

# covariate shift

## Covariate Shift or Feature Bias

- However, no chance for generalization if training and test samples have nothing in common.

$$P_{train}(\boldsymbol{x}, y) \neq P_{test}(\boldsymbol{x}, y)$$

- Covariate shift:
  - Input distribution changes
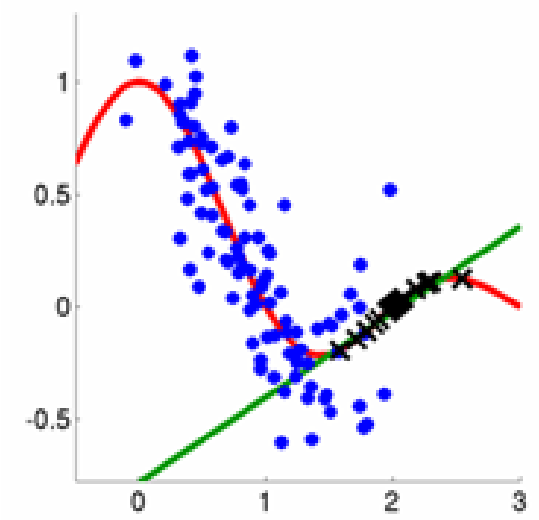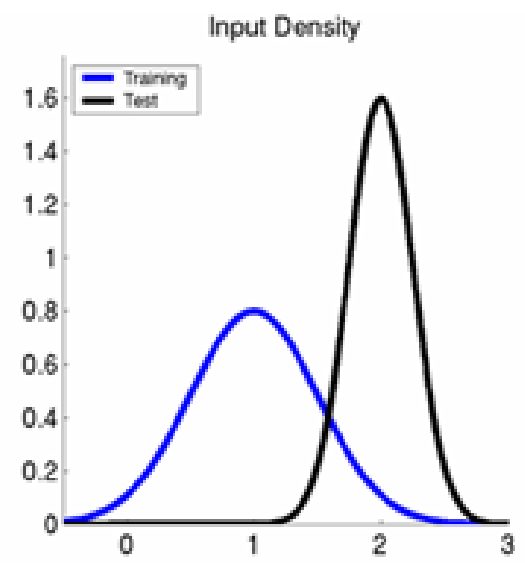  $$P_{train}(\boldsymbol{x}) \neq P_{test}(\boldsymbol{x})$$
  - Functional relation remains unchanged
  $$P_{train}(y|\boldsymbol{x}) = P_{test}(y|\boldsymbol{x})$$

# Covariate Shift

Training and test input follow different distributions, but functional relation remains unchanged.

Target Function $f(x)$
Learned Function $\widehat{f}(x)$
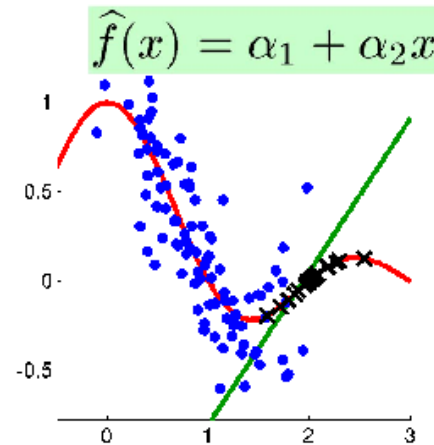● Training Sample $(x_i, y_i)$
✕ Test Sample $(t_i, u_i)$



Input Density

Goal: Estimate test output from $\{(x_i, y_i)\}_{i=1}^{n}$

# Importance-Weighted Least-Squares

$$\min_{\boldsymbol{\alpha}} \left[ \sum_{i=1}^{n} \frac{p_{test}(\boldsymbol{x}_i)}{p_{train}(\boldsymbol{x}_i)} \left( \widehat{f}(\boldsymbol{x}_i) - y_i \right)^2 \right]$$

- ■ IWLS is consistent even under covariate shift.

- ■ The idea is applicable to any likelihood-based methods!

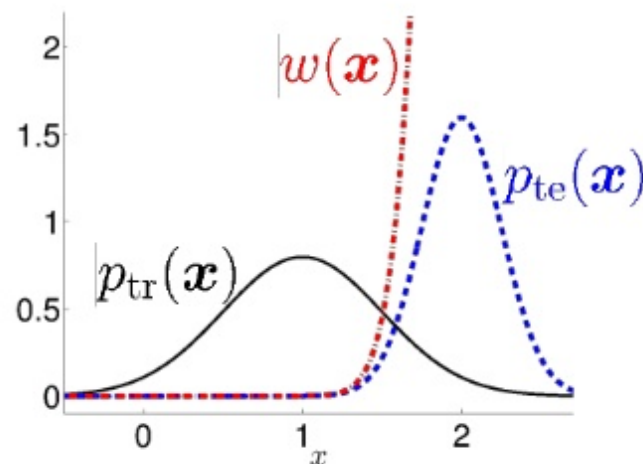  - ● Support vector machine, logistic regression, conditional random field, etc.



$\widehat{f}(x) = \alpha_1 + \alpha_2 x$

# A Problem in Covariate Shift Adaptation

■ Importance weight

$$w(\boldsymbol{x}) = \frac{p_{\text{te}}(\boldsymbol{x})}{p_{\text{tr}}(\boldsymbol{x})}$$

can **diverge to infinity** under a rather simple setting.

Cortes et al. (NIPS 2010)



$$p_{\text{tr}}(\boldsymbol{x}) = N(1, 0.5^2)$$
$$p_{\text{te}}(\boldsymbol{x}) = N(2, 0.25^2)$$

In this situation, the covariate shift adaptation is unstable since **estimated importance weight is unstable**☹