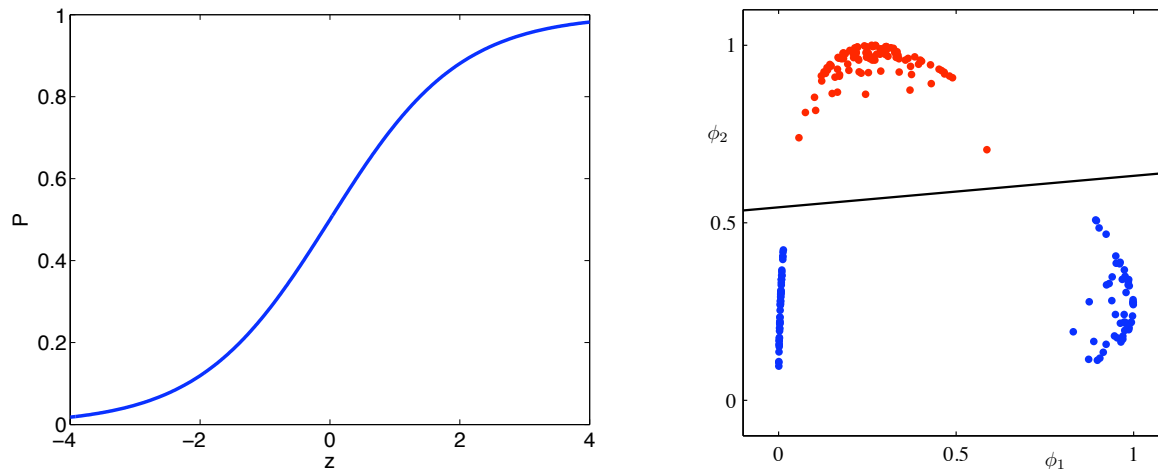


# Review

- Linear separability (and use of features)
- Class probabilities for linear discriminants
  - ▶ sigmoid (logistic) function
- Applications: USPS, fMRI



# Review

NB

$P(X, Y)$

$P(Y | X)$

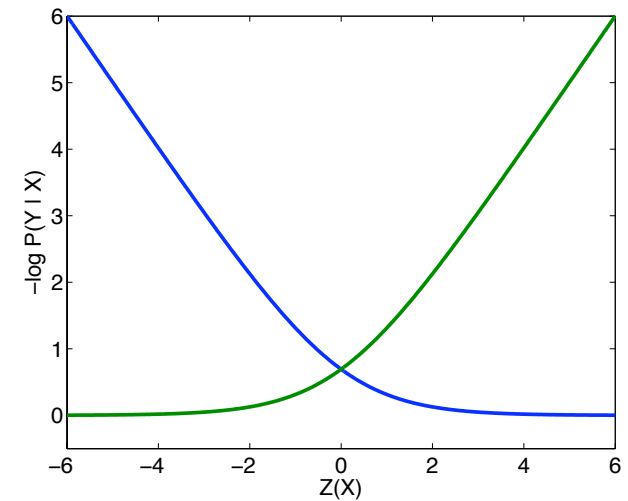
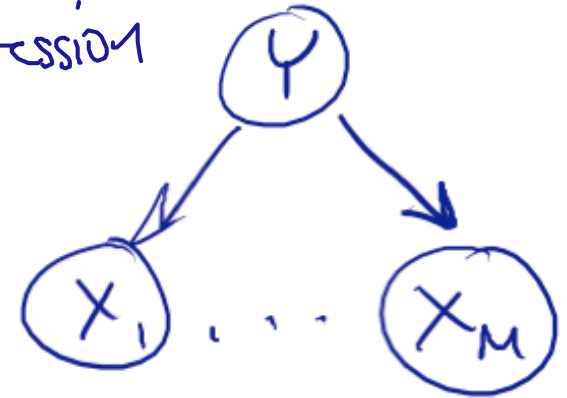
logistic regression

- Generative vs. discriminative
  - ▶ maximum conditional likelihood

- Logistic regression

- Weight space

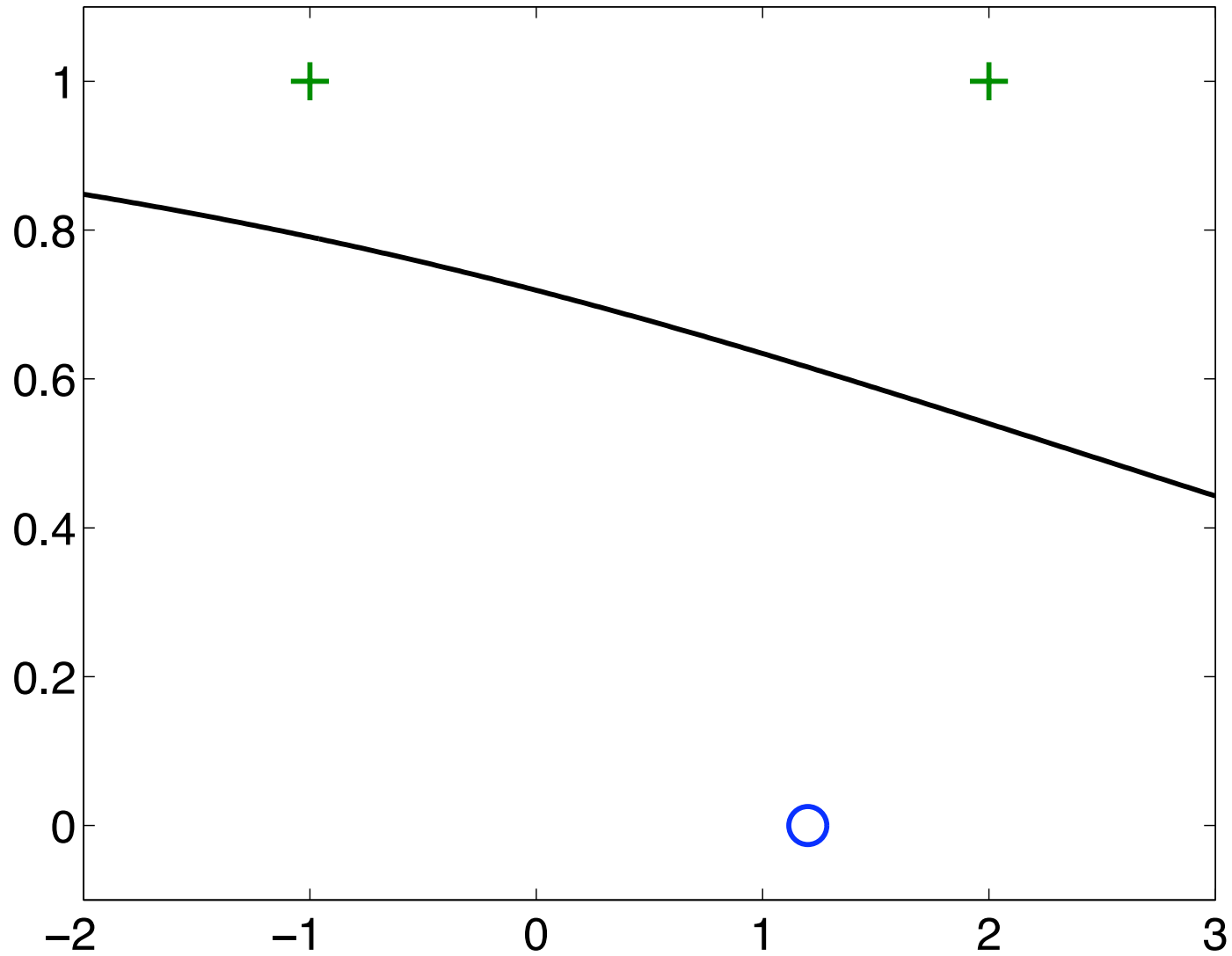
- ▶ each example adds a penalty to all weight vectors that misclassify it
- ▶ penalty is approximately piecewise linear



$$\approx \max(0, z) \quad \max(0, -z)$$

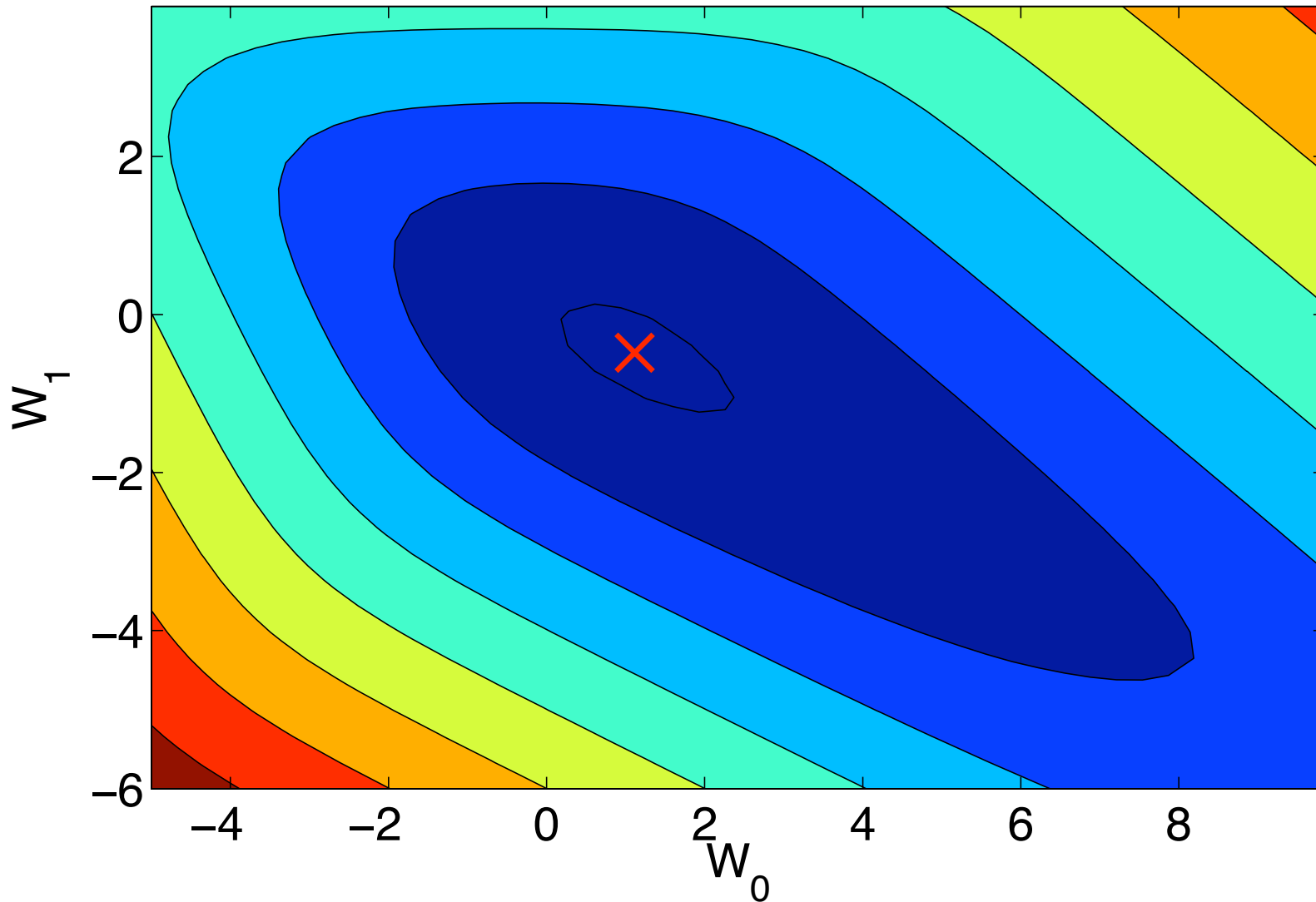
$$\approx \max(0, z) \quad \max(0, -z)$$

# Example



$$= \sum_{i=1}^3 \log \sigma(y^i x^i \cdot w)$$

$$-\log(P(Y_{1:3} | X_{1:3}, W))$$



# Generalization: multiple classes

$$k = \arg \max_{k'} z_k \geq z_{k'} \quad \forall k'$$

$$w^1 \dots w^C \in \mathbb{R}^{n+1}$$

- One weight vector per class:  $Y \in \{1, 2, \dots, C\}$

$$\triangleright P(Y=k) = \frac{\exp(z_k)}{\sum_{k'=1}^C \exp(z_{k'})}$$

$$\triangleright Z_k = w_0^k + \sum_{j=1}^n w_j^k X_j$$

$$\hat{Y} = \arg \max_{Y \in \{1, \dots, C\}} P(Y|X)$$

- In 2-class case:

$$\hat{Y} = \arg \max_{Y \in \{1, \dots, C\}} P(Y|X, w)$$

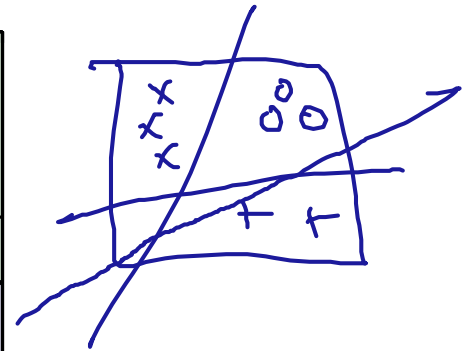
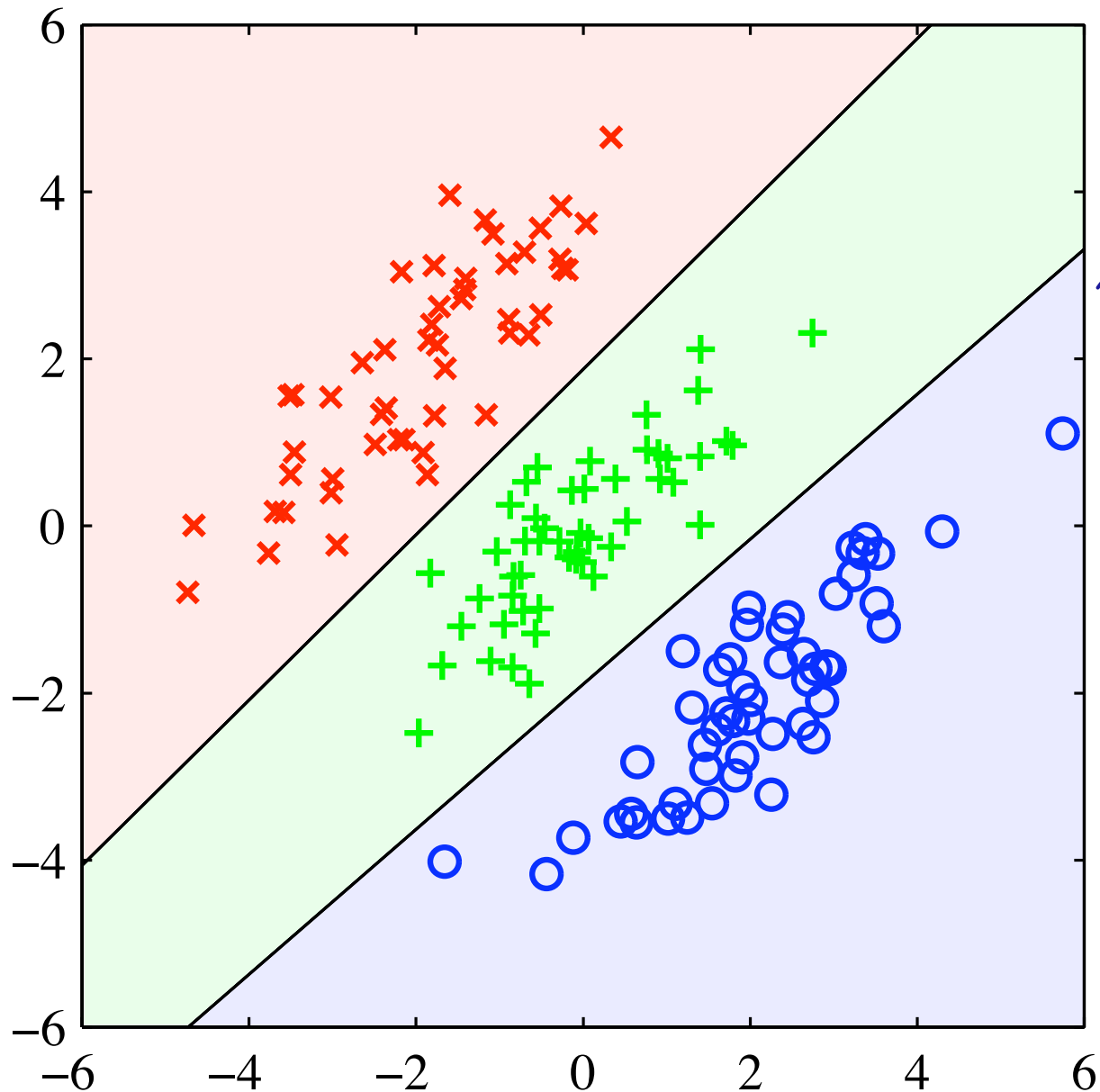
$$\hat{Y} = k \text{ iff } z_k \geq z_{k'} \quad (\forall k')$$

$$\frac{\exp(z_0)}{\exp(z_0) + \exp(z_1)}$$

$$= \frac{1}{1 + \exp(z_1 - z_0)}$$

$$z_1 - z_0 = \underbrace{w_0^1 - w_0^0}_{-w_0} + \sum_j \underbrace{(w_j^1 - w_j^0)}_{-w_j} X_j$$

# Multiclass example



See reading  
for more  
examples

figure from book

# Priors and conditional MAP

(2 class)

- $P(Y | X, W) = \sigma(Yz)$

- ▶  $Z = w_0 + \sum_j w_j x_j$

- As in linear regression, can put prior on W

- ▶ common priors:  $L_2$  (ridge),  $L_1$  (sparsity)

~~$-\log P(w) \propto \|w\|_2^2$~~   $\uparrow$

~~$-\log P(w) \propto \|w\|_1$~~

- $\arg \max_w P(W=w | X, Y) = \arg \max_w P(Y | X, W=w) P(W=w)$

→ actually  $-\log P(w) = k_1 + k_2 \|w\|_2^2$

or  $-\log P(w) = k_1 + k_2 \|w\|_1$

(linear r.t. affine function)

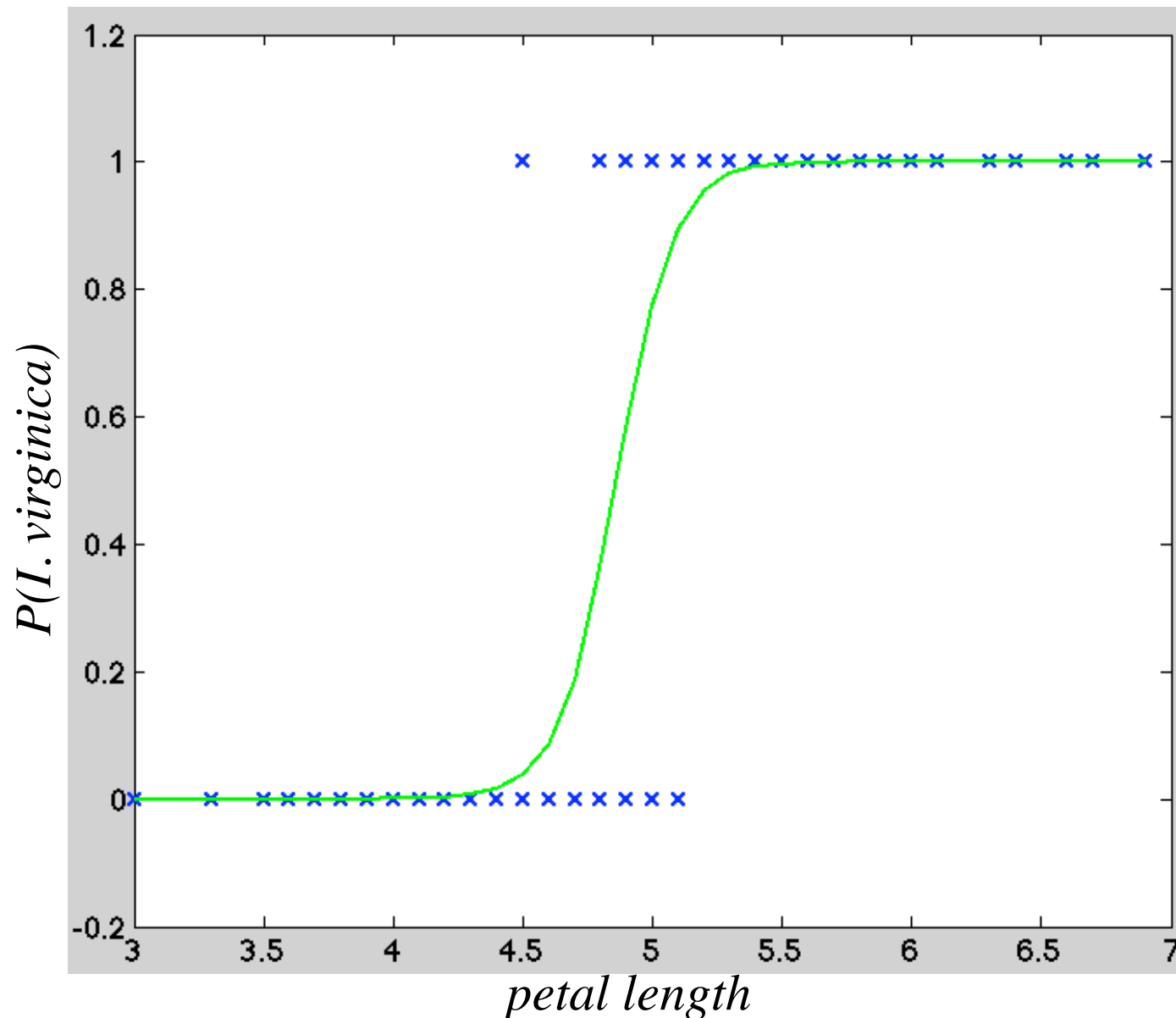
# Software

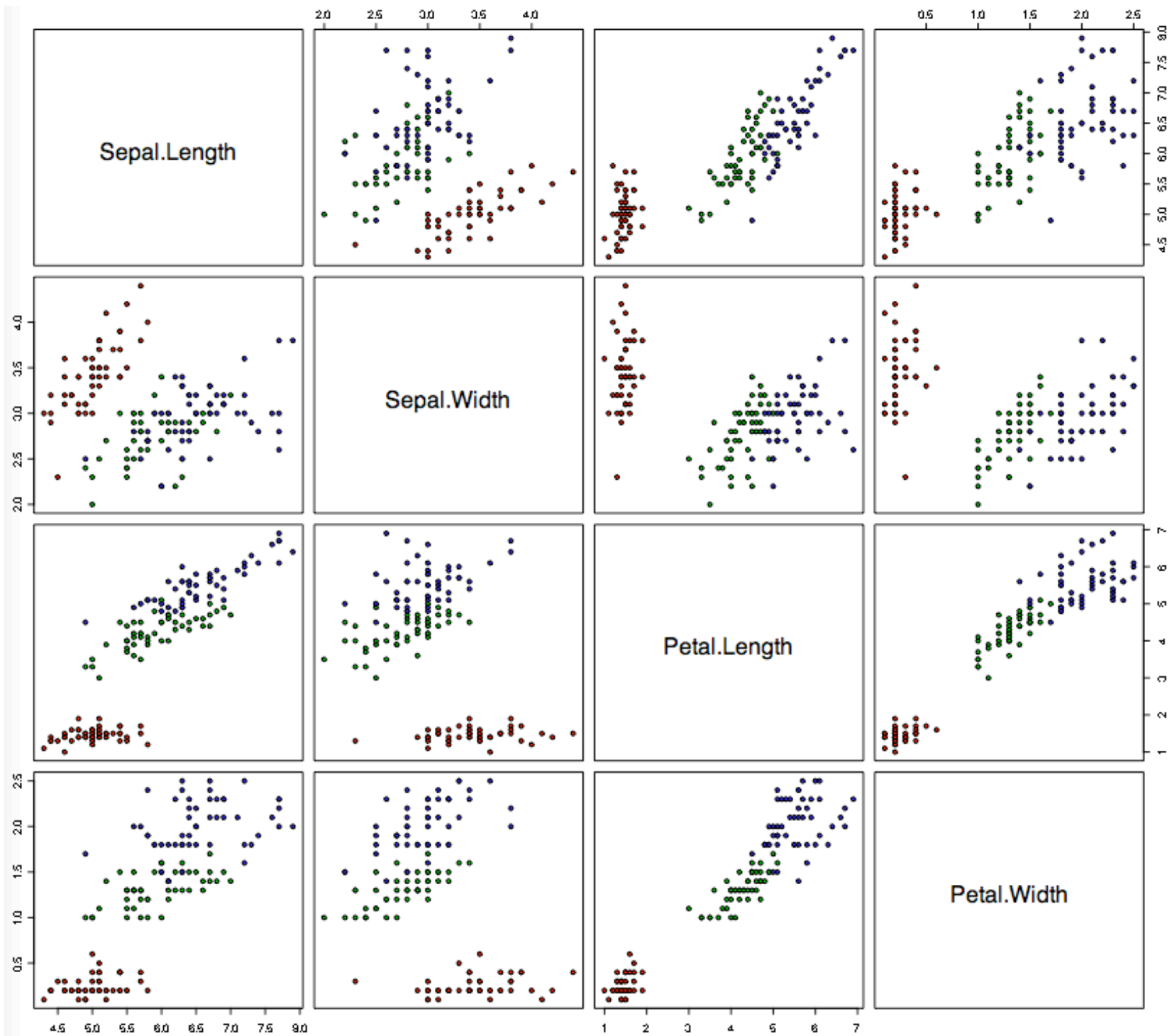
- Logistic regression software is easily available: most stats packages provide it
  - ▶ e.g., `glm` function in R
  - ▶ or, <http://www.cs.cmu.edu/~ggordon/IRLS-example/>
- Most common algorithm: Newton's method on log-likelihood (or  $L_2$ -penalized version)
  - ▶ called “iteratively reweighted least squares”
  - ▶ for  $L_1$ , slightly harder (less software available)

↑ still convex



# Historical application: Fisher iris data

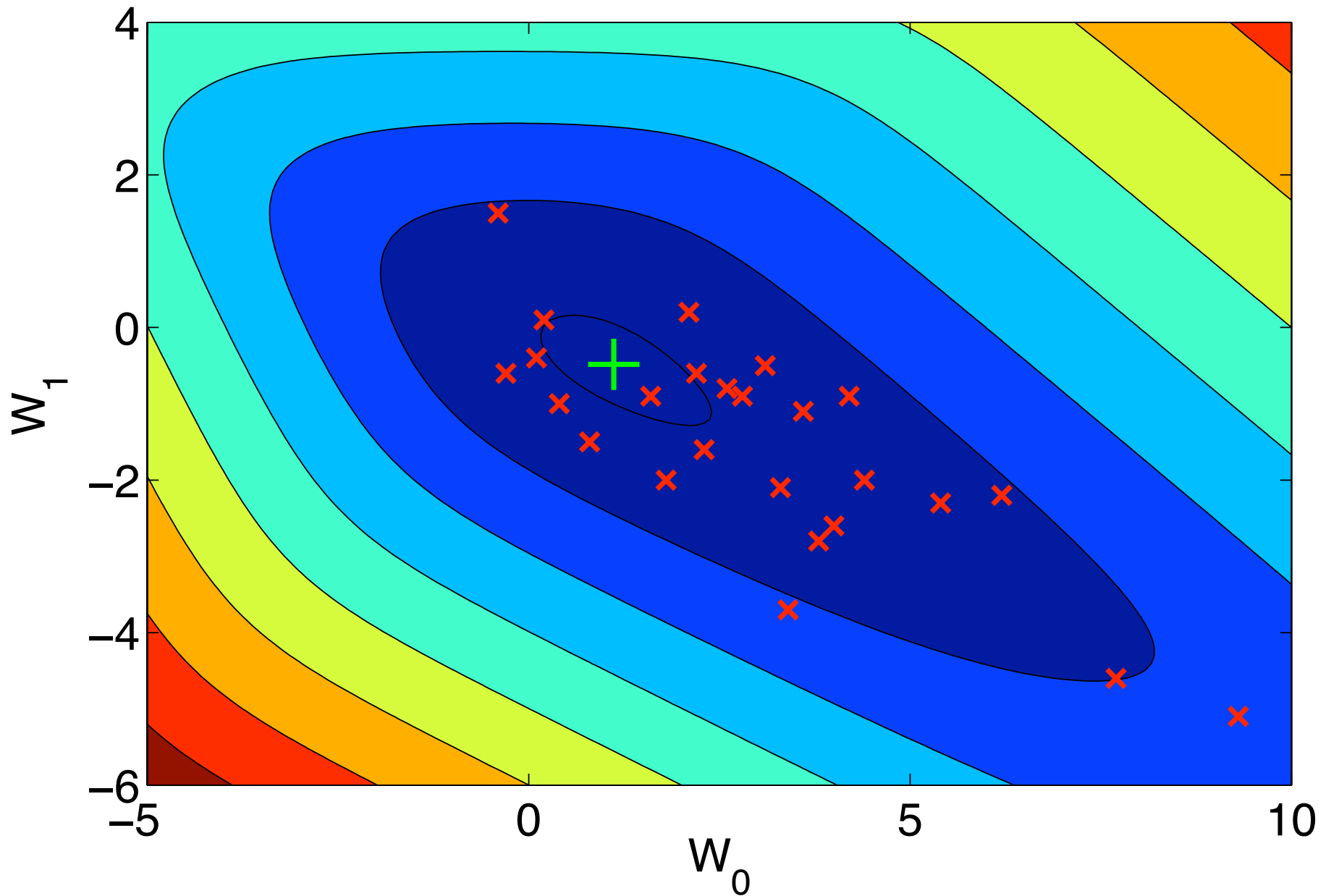




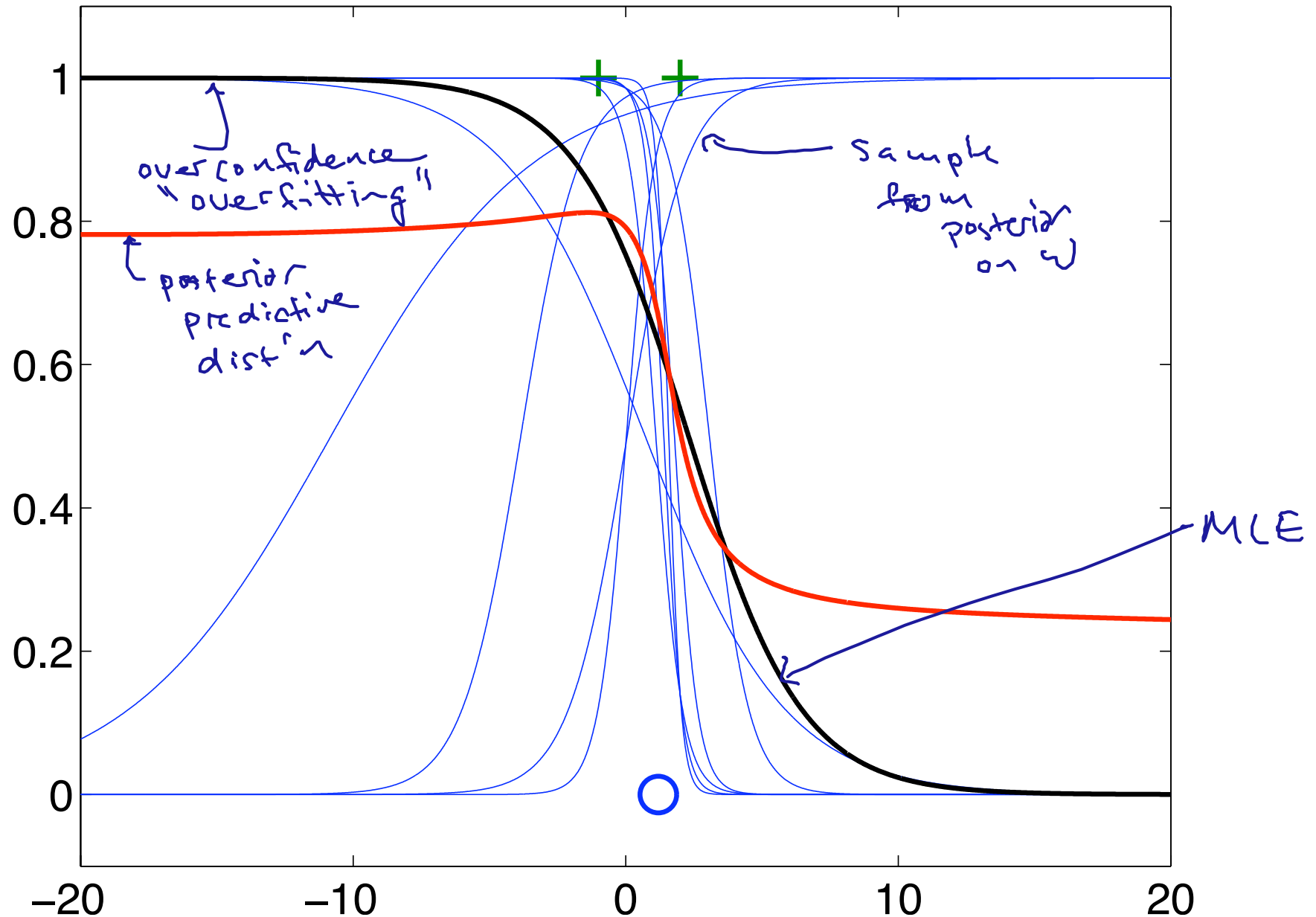
# Bayesian regression

- In linear and logistic regression, we've looked at
  - ▶ conditional MLE:  $\max_w P(Y | X, w)$
  - ▶ conditional MAP:  $\max_w P(W=w | X, Y)$
- But of course, a true Bayesian would turn up nose at both
  - ▶ why?

# Sample from posterior



# Predictive distribution



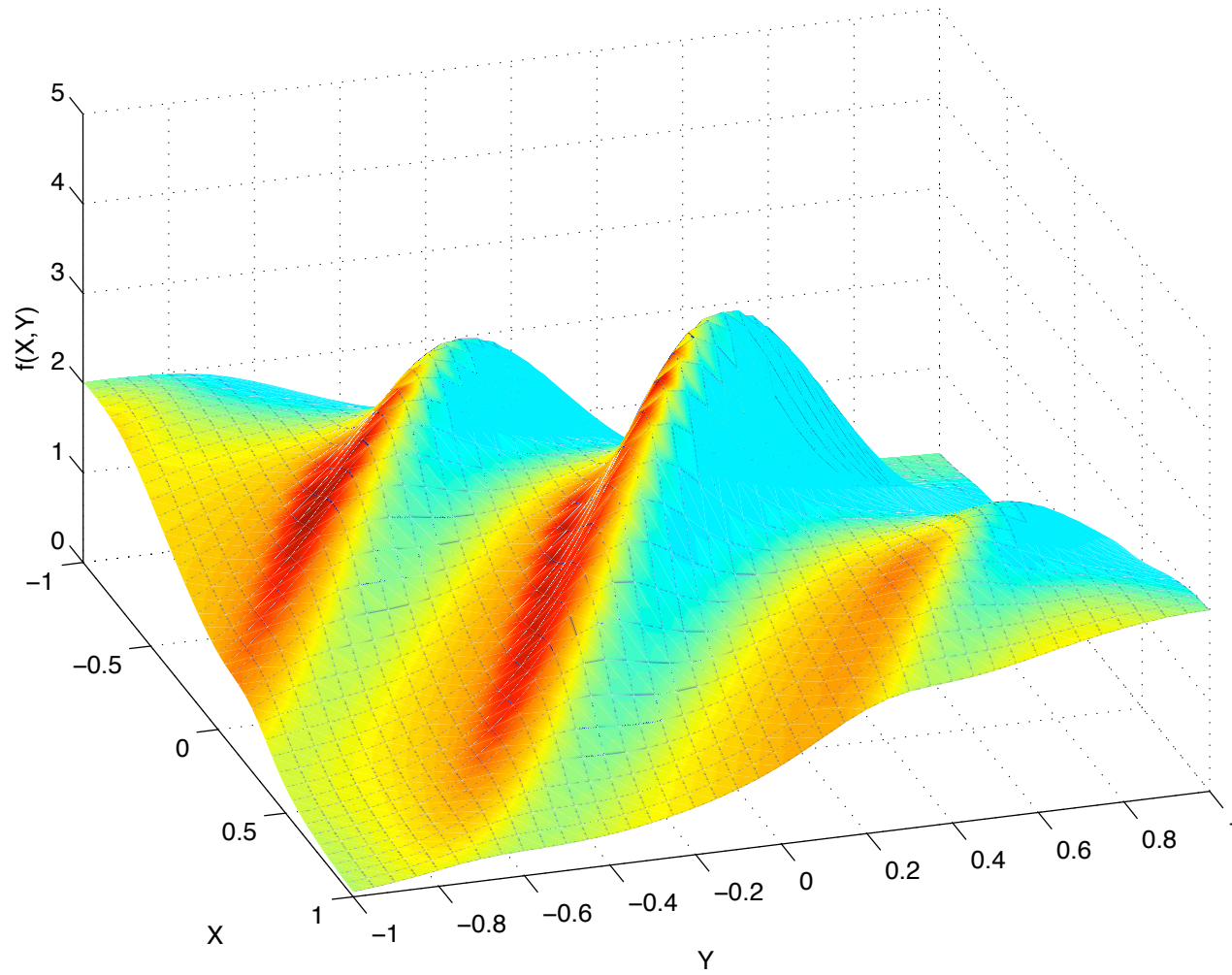
# Overfitting

- Overfit: training likelihood  $\gg$  test likelihood
  - ▶ often a result of overconfidence
- Overfitting is an indicator that the MLE or MAP approximation is a bad one
- Bayesian inference rarely overfits
  - ▶ may still lead to bad results for other reasons!
  - ▶ e.g., not enough data, bad model class, ...

# So, we want the predictive distribution

- Most of the time...
  - ▶ Graphical model is big and highly connected
  - ▶ Variables are high-arity or continuous
- Can't afford exact inference
  - ▶ Inference reduces to numerical integration (and/or summation)
  - ▶ We'll look at randomized algorithms

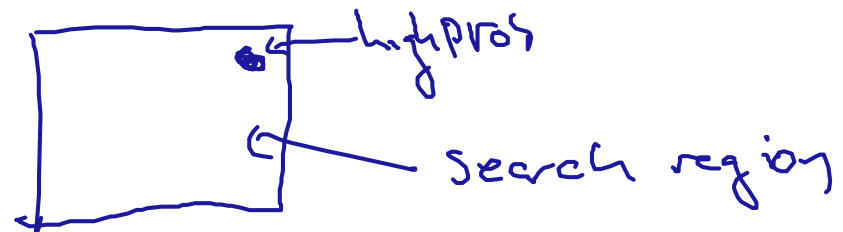
# Numerical integration



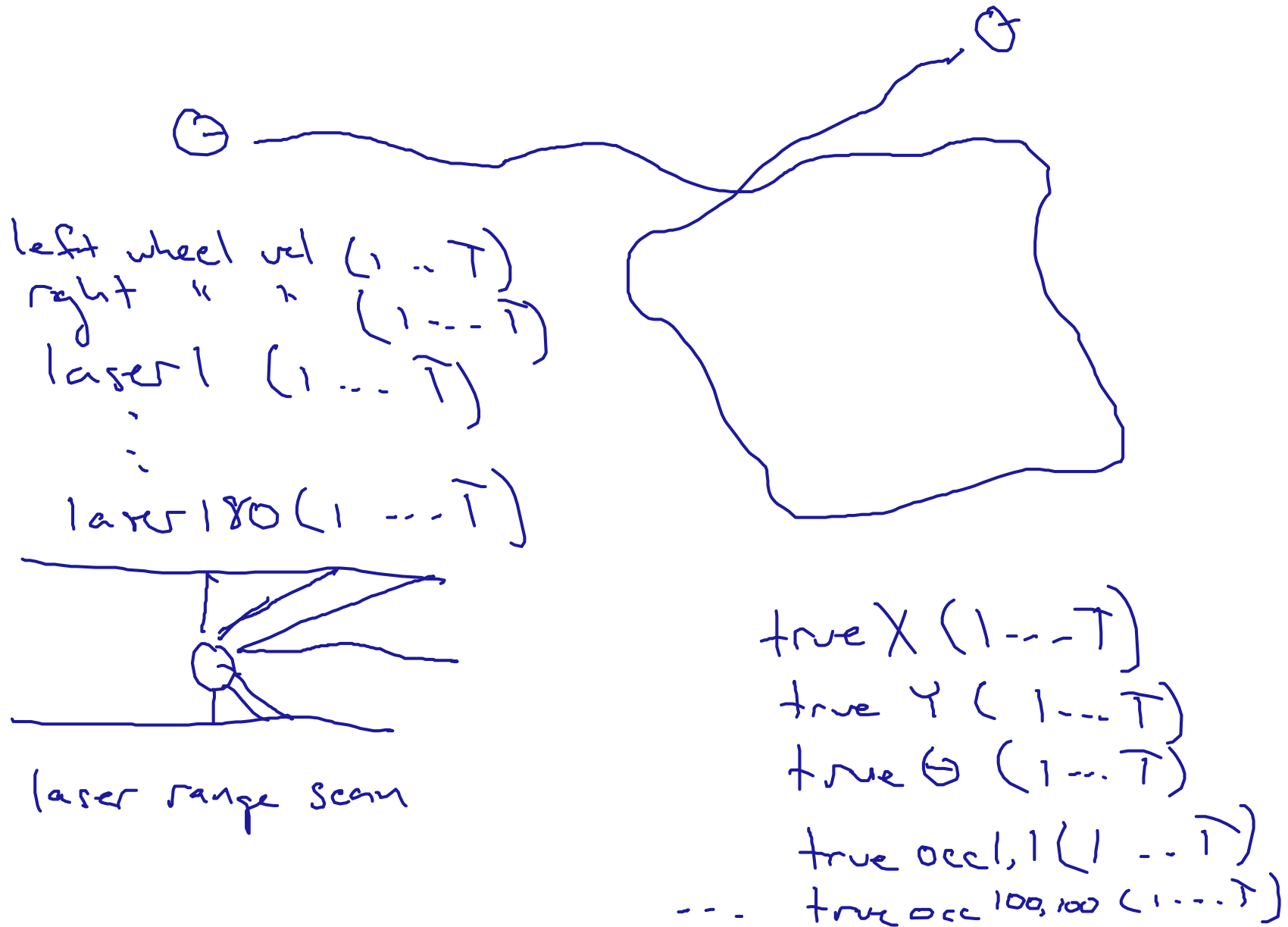


# 2D is 2 easy!

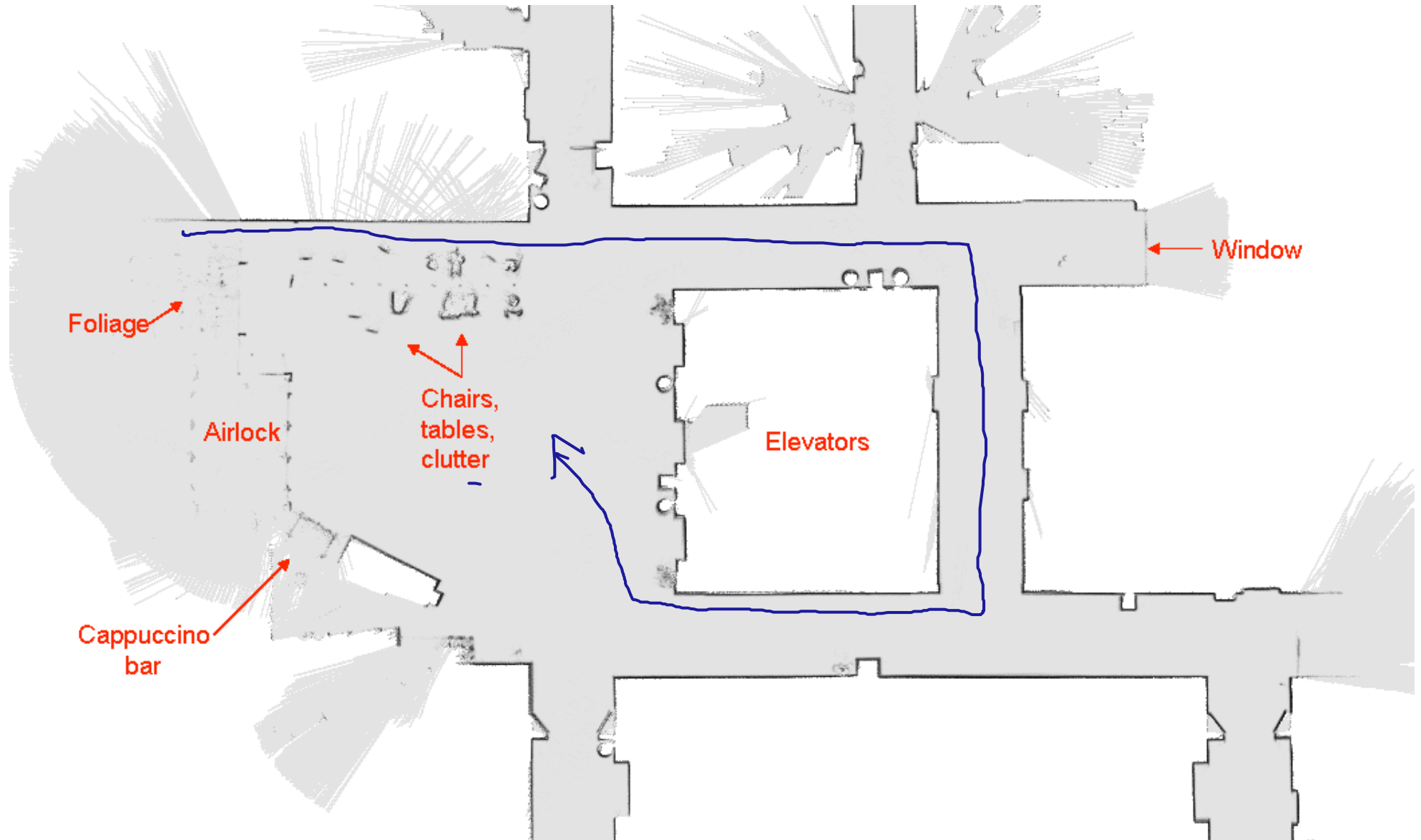
- We care about high-D problems
- Often, much of the mass is hidden in a tiny fraction of the volume
  - ▶ simultaneously try to discover it and estimate amount



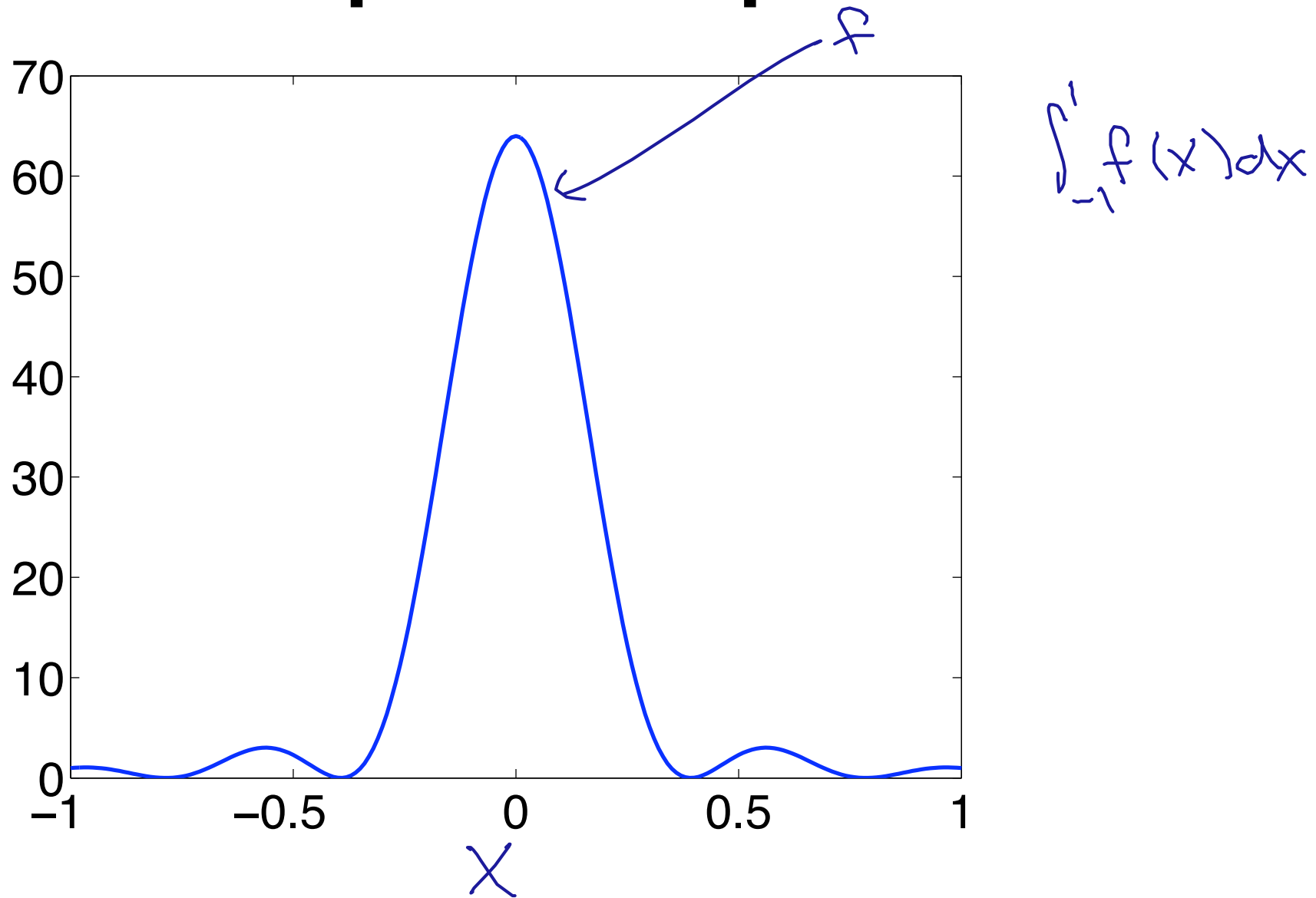
# Application: SLAM



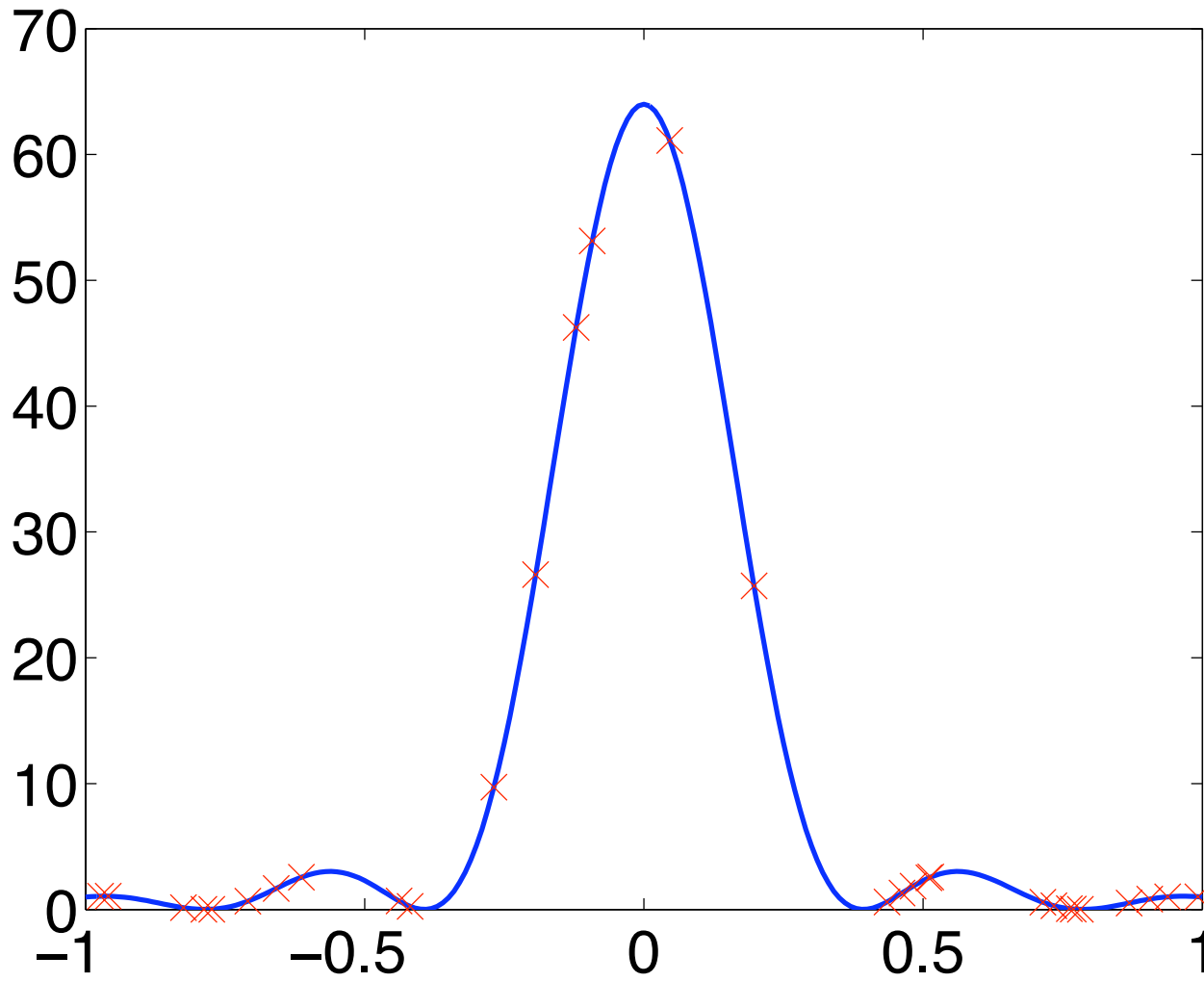
# Integrals in multi-million-D



# Simple 1D problem



# Uniform sampling



$$\int f dx \approx \hat{F} = \frac{1}{N} \cdot 2 \sum_{i=1}^N f(x_i)$$

$$\int f dx \approx 24.0$$
$$\hat{F} = 14.7$$

$$\int_{-1}^1 \frac{1}{2} dx = 1$$

# Uniform sampling

$$E(f(X)) = \int_{-1}^1 P(x) f(x) dx = \frac{1}{2} \int_{-1}^1 f(x) dx$$

↙  $\frac{1}{2} = \frac{1}{V} = \frac{1}{\text{volume}}$

- So,  $2 E(f(X))$  is desired integral (or,  $V E(f(x))$ )
- But standard deviation can be big
- Can reduce it by averaging many samples
- But only at rate  $1/\sqrt{N}$

# Importance sampling

- Instead of  $X \sim \text{uniform}$ , use  $X \sim Q(x)$
- $Q = \text{importance dist'n}$
- Should have  $Q(x)$  large where  $f(x)$  is large
- Problem: naive estimate  $\frac{1}{N} \sum_i f(X_i)$  wrong!  
 $E(f(X_i)) = \int Q(x) f(x) dx \neq \text{any simple fn of } \int f dx$

# Importance sampling

- Instead of  $X \sim \text{uniform}$ , use  $X \sim Q(x)$
- $Q =$
- Should have  $Q(x)$  large where  $f(x)$  is large
- Problem:

$$E_Q(f(X)) = \int Q(x) f(x) dx$$



# Importance sampling

$$h(x) \equiv f(x)/Q(x)$$

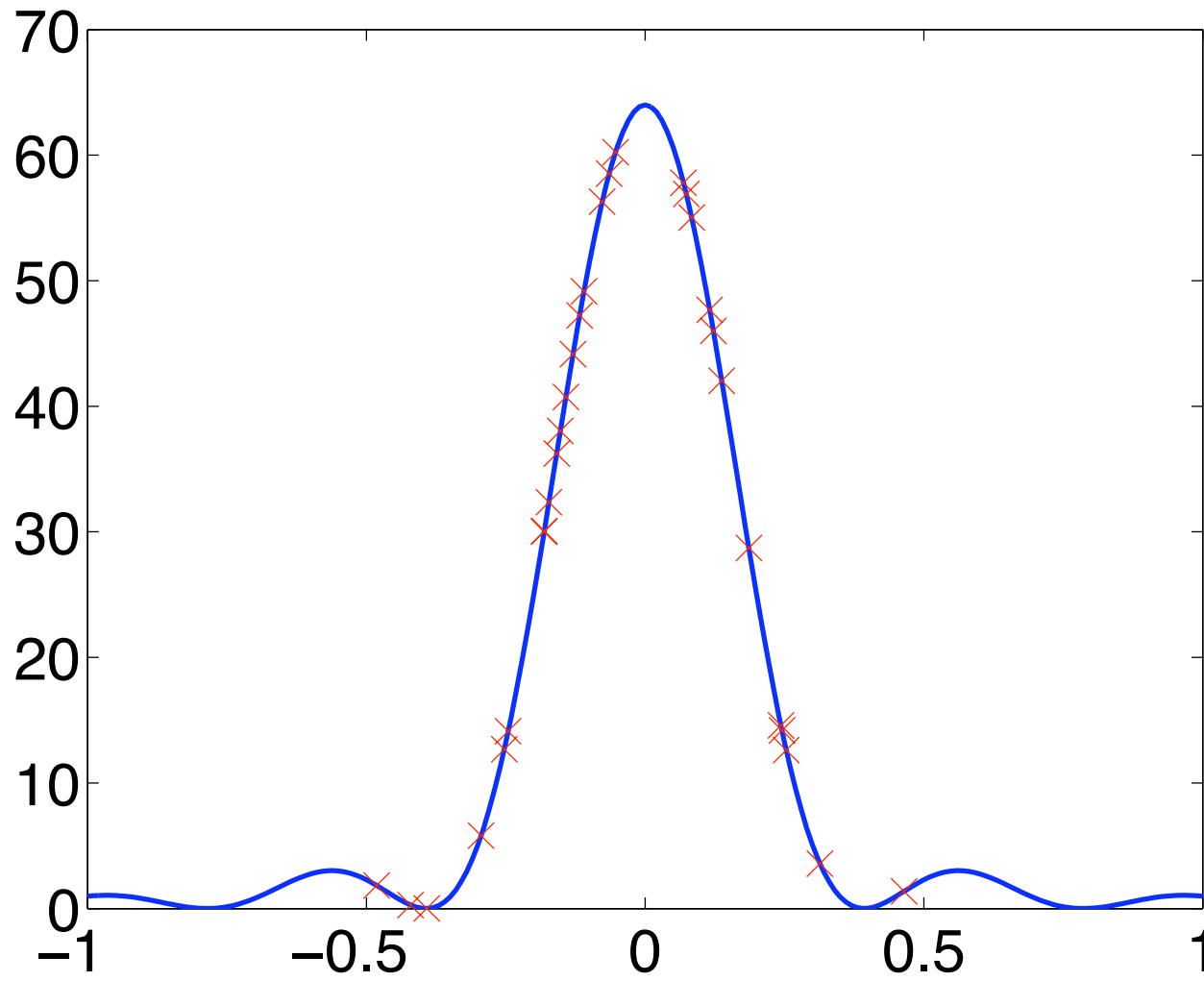
$$\begin{aligned} E_Q(h(x)) &= \int Q(x) h(x) dx \\ &= \int \cancel{Q(x)} \frac{f(x)}{\cancel{Q(x)}} dx \\ &= \int f(x) dx \end{aligned}$$

# Importance sampling

$$\hat{F} = \sum_{i=1}^N w(X_i) = \sum_{i=1}^N \underbrace{\frac{1}{Q(X_i)}}_{\text{importance weight}} f(X_i)$$

- So, take samples of  $h(X)$  instead of  $f(X)$
- $W_i = 1/Q(X_i)$  is **importance weight**
- $Q = 1/V$  yields uniform sampling

# Importance sampling



$$\int f dx = 24.0$$

$$\hat{F} = 25.8$$

$$Q(x) = N(x|0, .25^2)$$

30 samples

# Variance

- How does this help us control variance?
- Suppose:
  - ▶  $f$  big  $\Rightarrow Q$  big
  - ▶  $Q$  small  $\Rightarrow f$  small
- Then  $h = f/Q$ : moderate values
- Variance of each weighted sample is moderate
- Optimal  $Q$ ?  
$$\frac{f(x)}{Z} \leftarrow \text{never happens in practice}$$
$$= f(x) / \int f dx$$

# Importance sampling, part II

- Suppose we want

$$\int \underbrace{f(x)} \underbrace{dx = \int \underbrace{P(x)g(x)}_{\text{factors}} dx = E_P(g(X))}$$

- Pick N samples  $X_i$  from proposal  $Q(X)$
- Average  $W_i g(X_i)$ , where importance weight is
  - ▶  $W_i = P(x_i)/Q(x_i)$

# Importance sampling, part II

- Suppose we want

$$\int f(x)dx = \int P(x)g(x)dx = E_P(g(X))$$

- Pick N samples  $X_i$  from proposal  $Q(X)$
- Average  $W_i g(X_i)$ , where importance weight is
  - ▶  $W_i =$

$$E_Q(Wg(X)) = \int Q(x)[P(x)/Q(x)]g(x)dx = \int P(x)g(x)dx$$

# Two variants of IS

- Same algorithm, different terminology
  - ▶ want  $\int f(x) dx$  vs.  $E_P(f(X))$
  - ▶  $W = I/Q$  vs.  $W = P/Q$

# Parallel importance sampling

- Suppose we want

$$\int f(x)dx = \int P(x)g(x)dx = E_P(g(X))$$

- But  $P(x)$  is unnormalized (e.g., represented by a factor graph)—know only  $Z P(x)$

$P(x|\text{data})$   $\propto P(\text{data}|X) P(x)$   
sample from posterior



# Parallel IS

- Pick  $N$  samples  $X_i$  from proposal  $Q(X)$
- If we knew  $W_i = P(X_i)/Q(X_i)$ , could do IS

- Instead, set  $\hat{W}_i = Z P(X_i) / Q(X_i)$

► and,  $\bar{W} = \frac{1}{N} \sum_{i=1}^N \hat{W}_i$

- Then:  $E(\hat{W}_i) = \int \cancel{Q(x)} Z P(x) / \cancel{Q(x)} dx$   
 $= Z \int P(x) dx = Z$   
 $E(\bar{W}) = Z$  (lower variance)

# Parallel IS

- Final estimate:

$$\mathbb{E}_p(g(X)) \approx \hat{G} = \frac{1}{N} \sum_{i=1}^N \underbrace{\frac{\hat{w}_i}{\bar{w}}}_{\text{normalized importance wts}} g(x_i)$$