

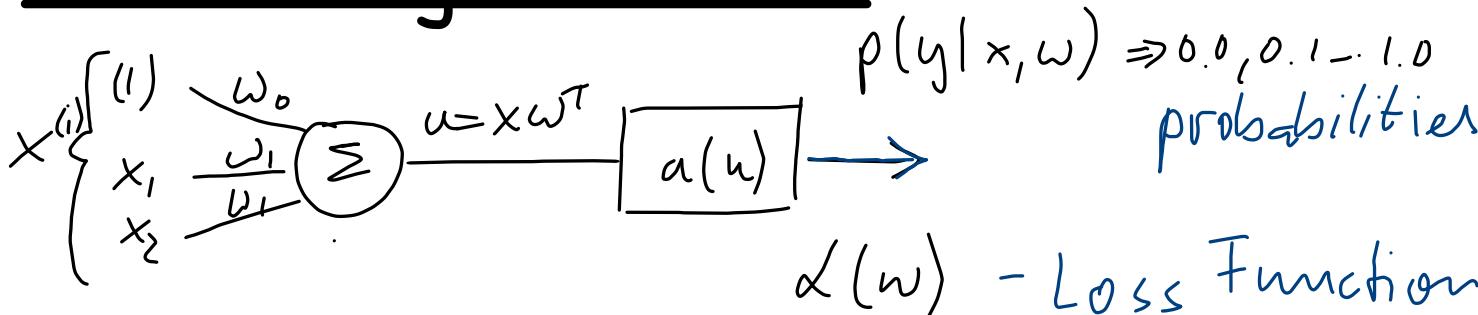
Bayesian Machine Learning
&

Its Application to Neural Networks

Questions

- Why Bayesian?
- Why do we need uncertainty estimates?
- What is overfitting? Where does it come from?
- What are we optimizing?
- Where does the Cost Function come from?
- How can we solve complicated Probability Distributions?

Motivating Example

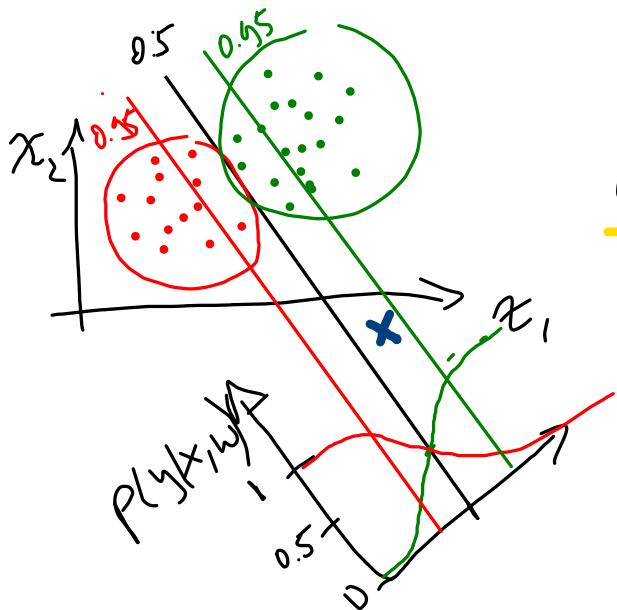


$\ell(w)$ - Loss Function

$\nabla_w \ell(w)$ - Gradient

$$\underline{w} \leftarrow w - \alpha \nabla_w \ell(w)$$

Gradient descent



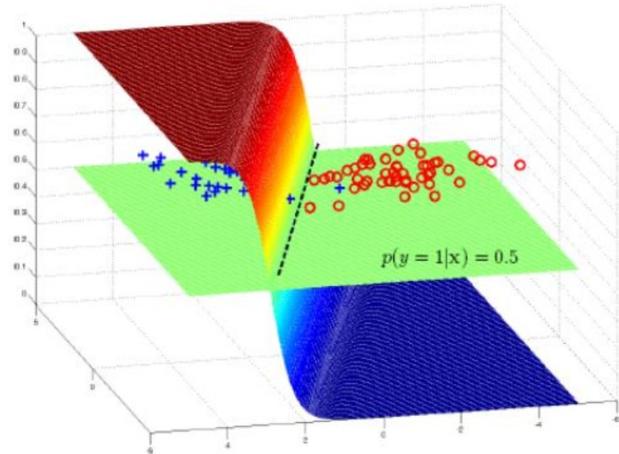
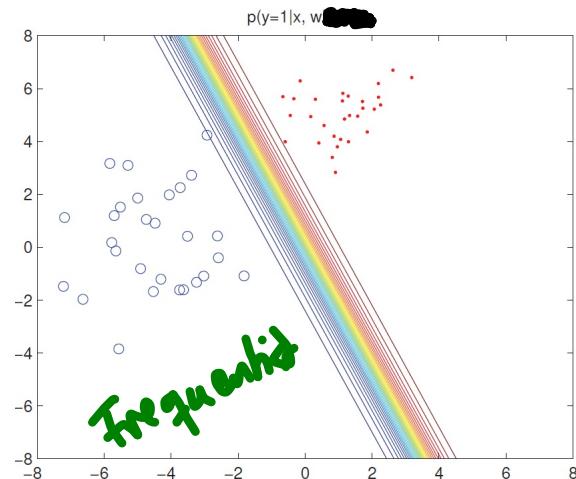
Frequentist
Point Estimation

Motivating Example

$$x^{(i)} \left\{ \begin{array}{l} x_0 \\ x_1 \\ x_2 \end{array} \right. \xrightarrow{\omega_0} u = x \omega^T \xrightarrow{\sum} a(u) \rightarrow p(y|x, \omega) \Rightarrow 0.0, 0.1 - 1.0 \text{ probabilities}$$

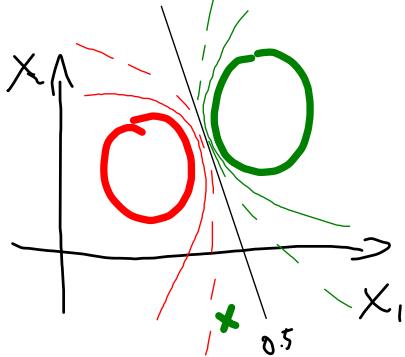
$\ell(w)$ - Loss Function

$\nabla_w \ell(w)$ - Gradient

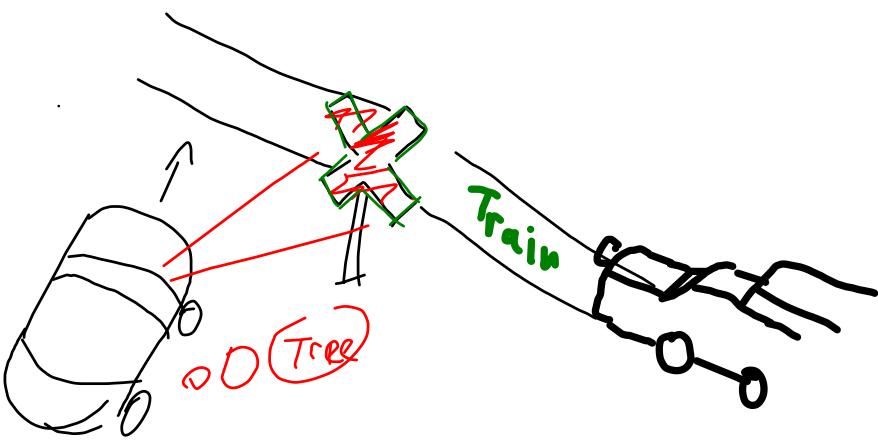
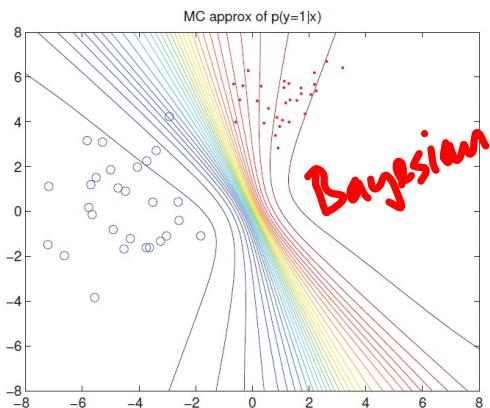


$\ell(w)$
t descent)

Uncertainty: Single Neuron \rightarrow Bayesian

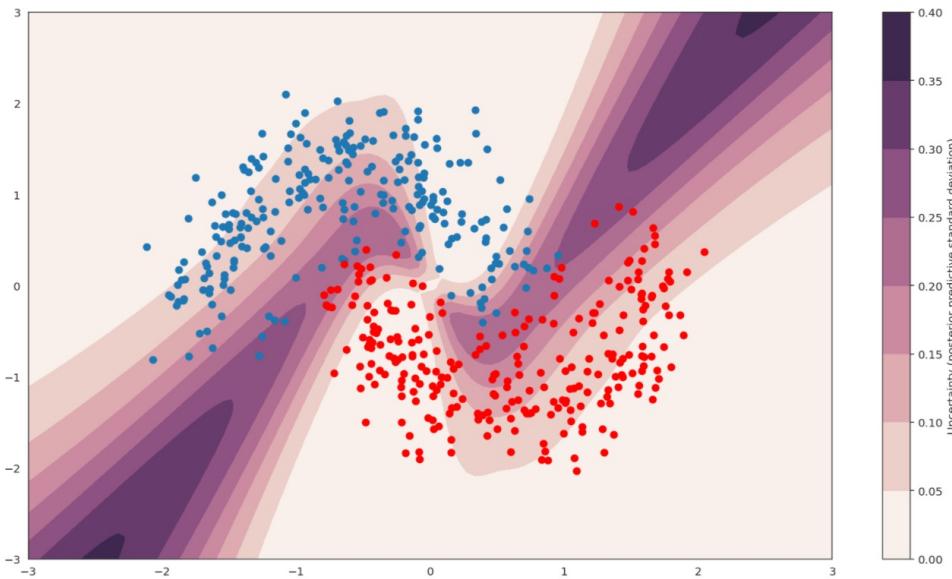
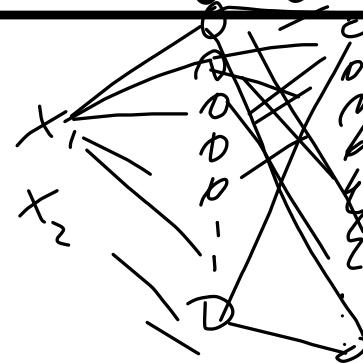


- Medical field
- Automotive

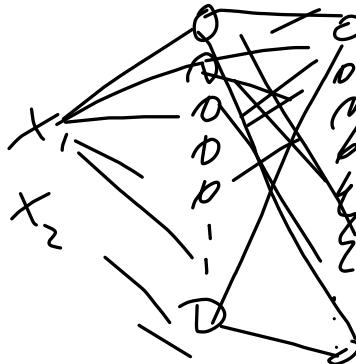
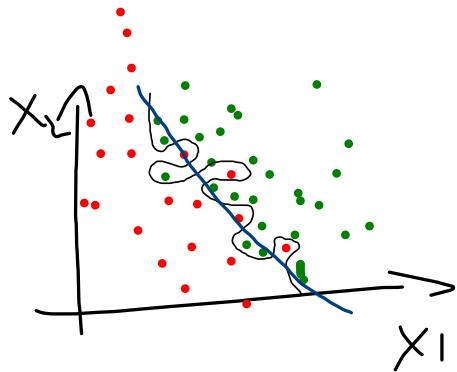


More Complex Data and Model

Neural Net



Over fitting



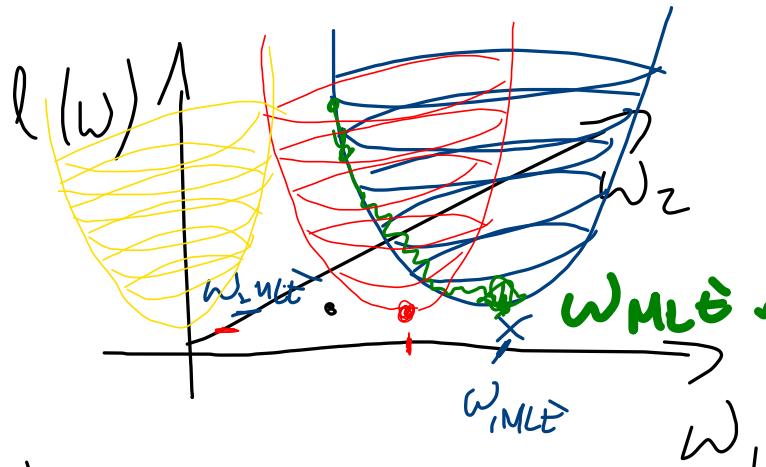
Optimization World

$$l_r(\omega) = l(\omega) + \underbrace{\alpha r(\omega)}_{\text{Regularizer}}$$

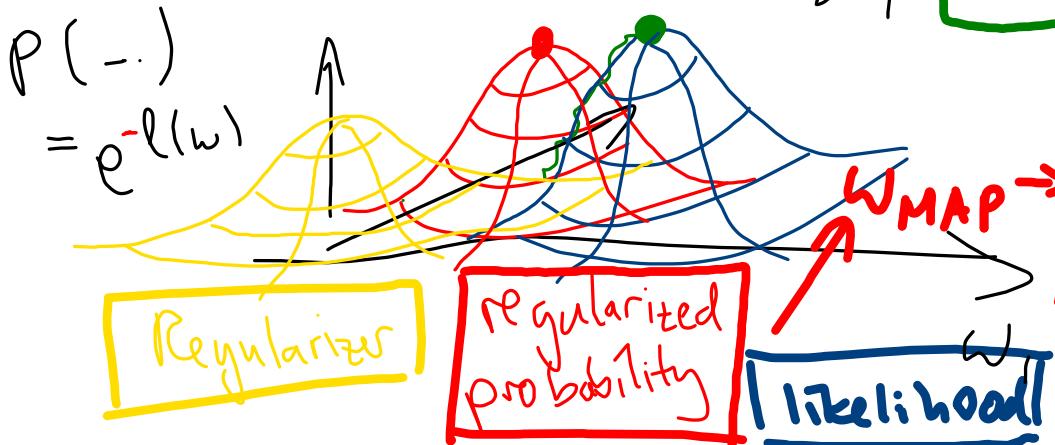
Regularizer

Bayesian Point Estimate

Point Estimation vs. Probability Distribution over Parameters (weights)



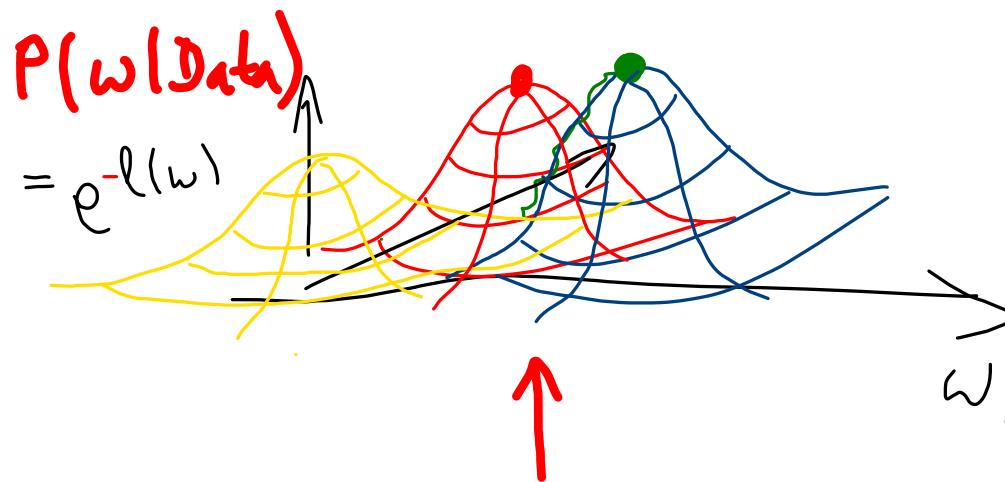
$w_{MLE} \rightsquigarrow$ Frequentist Point Estimate



$w_{MAP} \rightsquigarrow$ Bayesian Point Est..

Regularized Logistic Regression

Goal: Probability Distribution over Parameters



We want the whole distribution
over the weight,

instead of a point estimate

Bayesian Probability Intuition

Biased Coin Flips with „Small Data“

H T

↑
bend
coin

A

$y = [H H H H H H H H H H]$
 $\{1, 1, 1, 1, 1, 1, 1, 1, 1, 1\}$
 $n=10$

probability of getting Heads

$$\theta_H = \frac{10}{10} = 1$$

Is this correct?

Can it be 100%?

Black Swan Paradox

Biased Coin Flips with „Small Data”

A pair of handwritten symbols. The first symbol, on the left, looks like a capital letter 'H' with a diagonal line through it. The second symbol, on the right, looks like a capital letter 'A'.

$$y = \left[\begin{matrix} H & H & H & H & H & H & H & H & H & H \\ \{1,1,1,1,1,1,1,1,1,1\} \end{matrix} \right]_{1,10}$$

probability of getting Heads „Hallucinations“

$$\theta_H = \frac{10+1}{10+2} = \cancel{\underline{\underline{0.833}}}$$

Choose your „Hallucinations“ according to your „PRIOR“ knowledge!

Probabilistic Model of a Coin Flip

Overfitting with MLE

H 

$$P(y=\{ \dots \} | \theta) = \prod_{i=1}^n \text{Ber}(y^{(i)} | \theta) = \prod_{i=1}^n \theta^{y^{(i)}} \cdot (1-\theta)^{1-y^{(i)}}.$$

likelihood

$$= \theta^n \cdot (1-\theta)^{n_0}$$

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} P(y | \theta)$$

$$P(y | \theta=0.3) = \prod_{i=1}^{10} 0.3^1 = 0.3^{10} \times$$

$$P(y | \theta=0.5) = 0.5^{10} \times$$

$$P(y | \theta=1) = 1^{10} = 1 \quad \checkmark$$

Max

$y = [H, H, H, H, H, H, H, H, H, H]$
 $\{1, 1, 1, 1, 1, 1, 1, 1, 1, 1\}$

$$\underline{\theta^n \cdot (1-\theta)^{n_0}}$$

Short's
 $\hat{\theta}_{\text{MLE}}$ # 1's
flips

Black Swan Paradox

Overfitting - Bayesian Fix

H A

P($\theta | y$)

Posterior

Proportionality

$$\propto \prod_{i=1}^n \text{Ber}(y_i | \theta) \cdot "?"$$

$$\propto \theta^{n_1} \cdot (1-\theta)^{n_0} \cdot \theta^{\alpha-1} \cdot (1-\theta)^{\beta-1}$$

$$\boxed{\propto \theta^{n_1 + \alpha - 1} \cdot (1-\theta)^{n_0 + \beta - 1}}$$



normalization
constant $\frac{1}{Z}$
is missing

$$\Rightarrow \theta_{\text{Bayes}} = \frac{\# 1's + \alpha}{\# \text{flips} + \alpha + \beta}$$

$$[\alpha=1; \beta=1]$$

$$\theta_{\text{Bayes}} = \frac{10+1}{10+1+1} = \frac{10}{12} = \underline{\underline{0.833}}$$

Beta - Bernoulli Model

$P(\theta | y)$

Posterior

$$\propto \prod_{i=1}^n \text{Ber}(y_i | \theta)$$

$$\propto \theta^{n_1} \cdot (1-\theta)^{n_0}$$

"?"

$$\cdot \theta^{\alpha-1} \cdot (1-\theta)^{\beta-1}$$

$$\boxed{\propto \theta^{n_1+\alpha-1} \cdot (1-\theta)^{n_0+\beta-1}}$$

$$P(\theta | y) = \text{Ber}(y | \theta) \cdot \text{Beta}(\alpha, \beta)$$

$$= \text{Beta}(\alpha + n_1, \beta + n_0)$$

$$= \text{Beta}(\alpha', \beta')$$

} Conjugate Prior

Prior - Posterior

deterministic

$$\theta_i$$



$$\boxed{\arg \max_{\theta} \text{Ber}(y|\theta)}$$

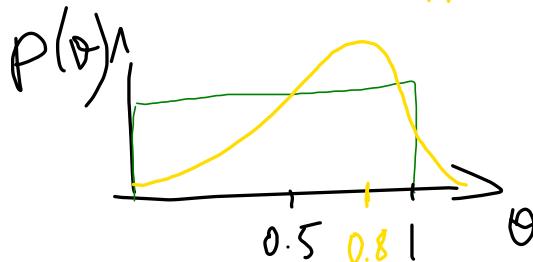
$$\theta_{MLE} = \delta(\theta_{MLE})$$

dirac

$$\theta \sim p(\theta)$$

random
Priors:

i.e. Beta(α, β)



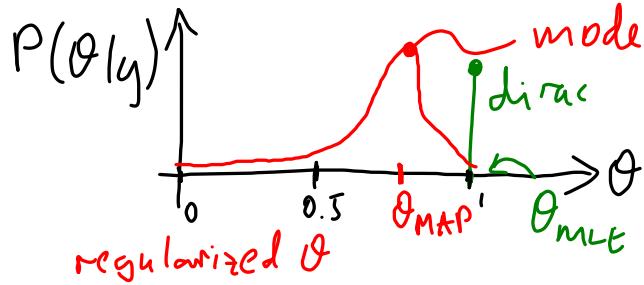
Like likelihood

Norm.

Posterior

$$\theta|y \sim p(\theta|y)$$

Beta($\alpha+n_1, \beta+n_0$)



regularized θ

mode
dirac
 θ_{MAP}
 θ_{MLE}

General Bayes Formula

$$p(\theta|y) = \frac{p(y|\theta) p(\theta)}{p(y)}$$

(Model) likelihood prior (Hallucinations)

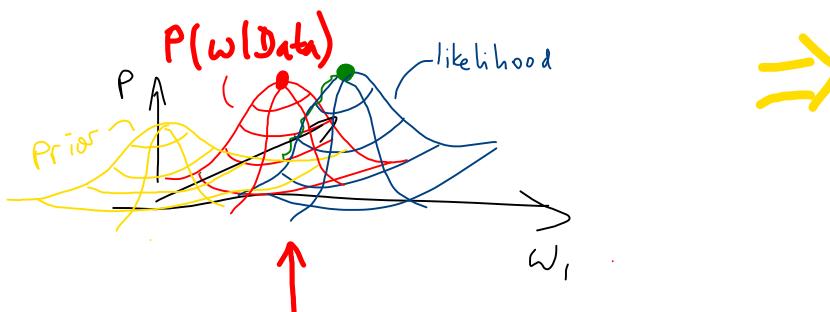
Posterior
over
parameters
(weights)

normalization
constant
aka. prob. of data

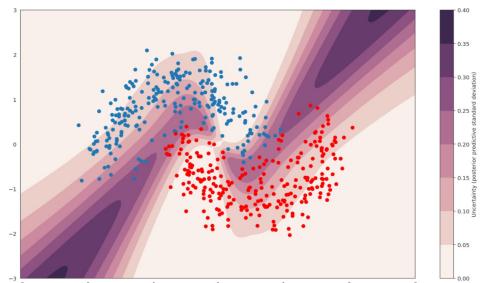
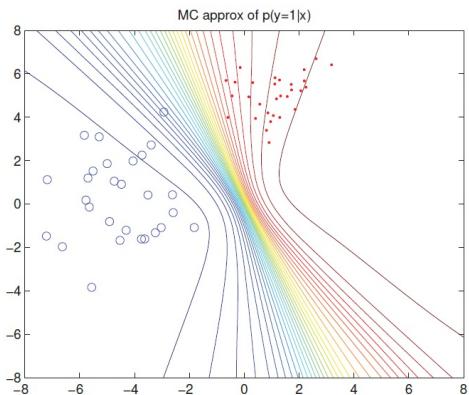
→ makes it sum to one

Remember?

Goal: Probability Distribution over Parameters



We want the whole distribution
over the weights,
instead of a point estimate



Easy Posterior vs. Nasty Posterior

Models with a Closed Form Solution
aka. Conjugate Prior

Binary Data • Beta-Bernoulli

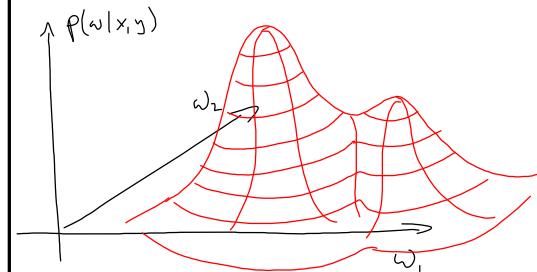
Multiclass Data • Dirichlet-Categorical

Linear Regression • Gaussian-Gaussian

•
•
•

Difficult/Intractable
Postiors

- Logistic Regression/
Single Neuron
- Neural Networks



How to solve Nasty Posterior Distributions

→ Approximate Methods:

- Variational Methods (i.e. VAE)
- Monte Carlo Methods

Part II

Connecting Bayes
to the Optimization
of a Single Newton

General Bayes Formula - Coin Model

$$p(\theta|y) = \frac{p(y|\theta) p(\theta)}{p(y)}$$

(Model) likelihood prior (Hallucinations)

Posterior over parameters (weights)

normalization constant
aka. prob. of data

→ makes it sum to one

General Bayes Formula - Single Neuron

$$p(\theta | x, y) = \frac{p(y | x, \theta) p(\theta)}{Z}$$

?

(Model) likelihood

?

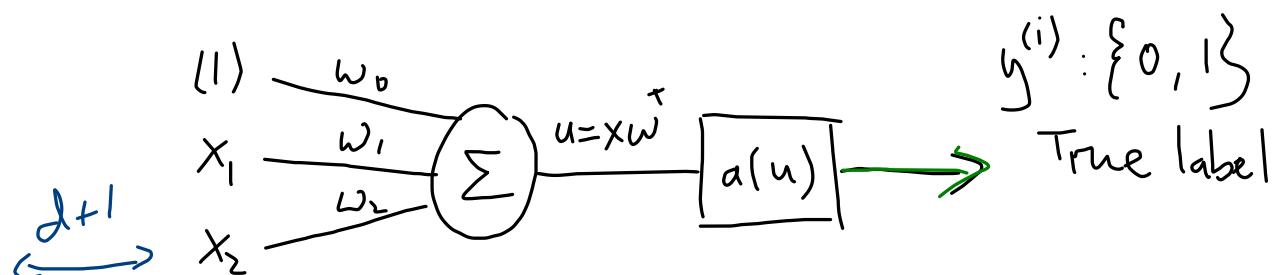
prior (Hallucinations)

↑

Also adding features x

↖ normalization const.
integral of numerator

Recap: Single Neuron Optimization



$$x = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ \vdots & \vdots & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} \end{bmatrix} \quad n$$

$$\omega = (w_0, w_1, w_2)$$

$$a(u) = \hat{y} \quad \text{Prediction}$$

Learning as Optimization

Regulated Loss Function

$$l_r(\omega) = e(\omega) + \alpha r(\omega)$$

Gradient
Update

$$\frac{\nabla_{\omega} l(\omega)}{w \leftarrow w - \alpha \nabla_{\omega} l(\omega)}$$

From Regularized Loss Function to Bayes

$$l_r(\omega) = l(\omega) + \alpha r(\omega) \quad | \exp(...)$$

$$\frac{-l_r(\omega)}{e} = \boxed{e^{-l(\omega)} \cdot e^{-\alpha r(\omega)}}$$

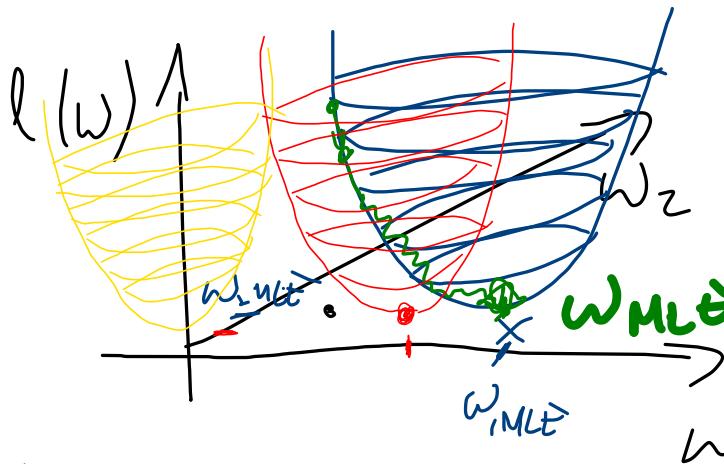
(Model) likelihood prior (Hallucinations)

$$p(\theta|x,y) = \frac{p(y|x,\theta) p(\theta)}{Z}$$

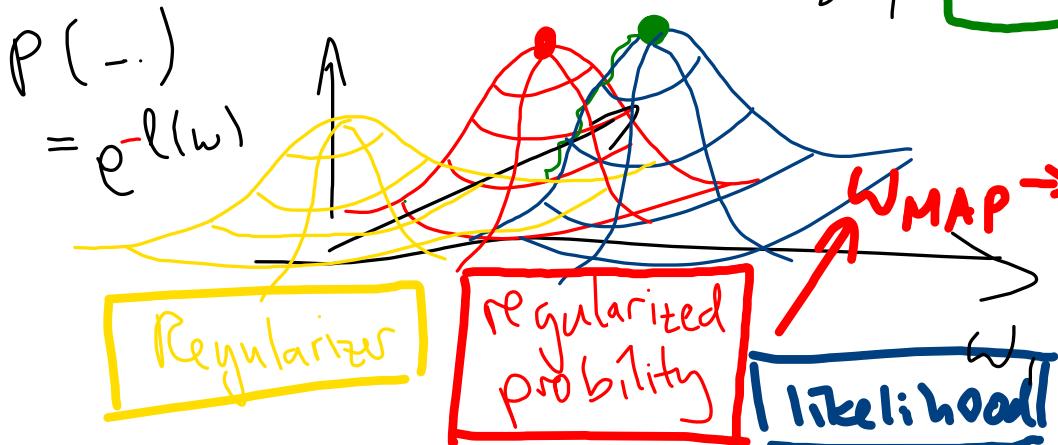
Z ↗ normalization const.
integral of numerator

$$\int_{-\infty}^{\infty} e^{-l(\omega)} \cdot e^{-\alpha r(\omega)} d\omega$$

Remember this Connection?



$w_{MLE} \rightsquigarrow$ Frequentist
Point Estimate

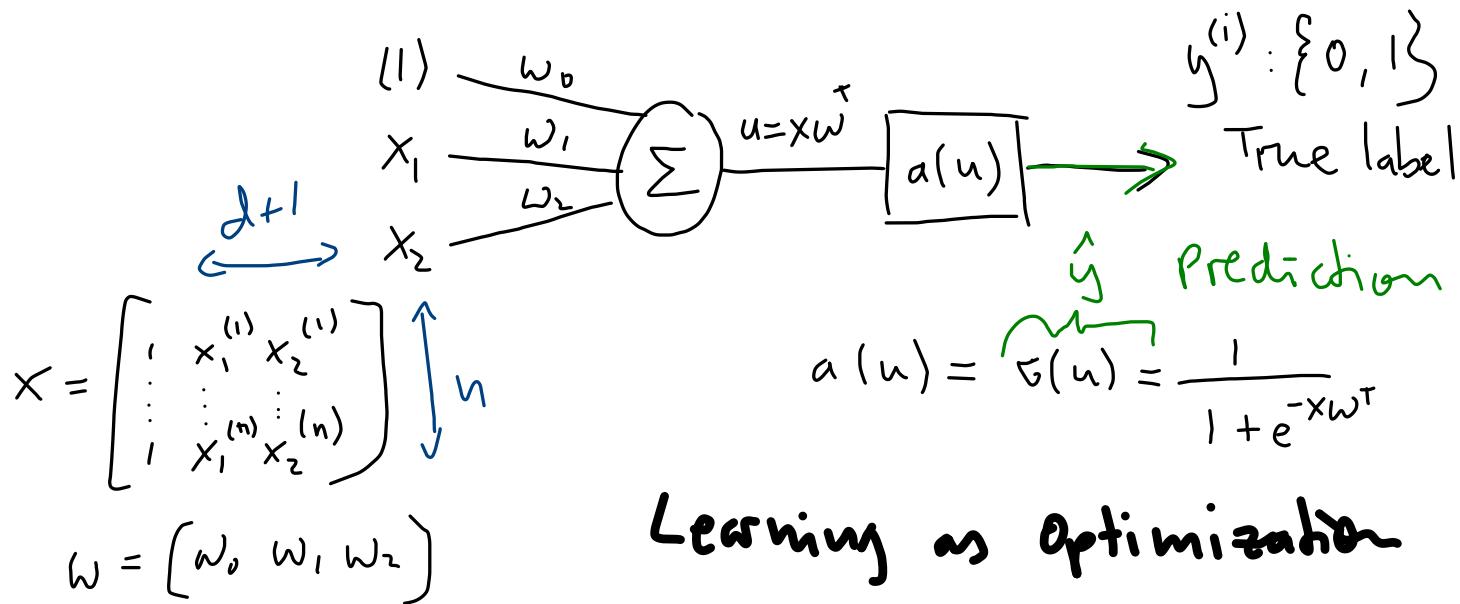


$w_{MAP} \rightsquigarrow$ Bayesian
Point Est.

Regularized
Logistic Regression

Going Deeper
into the Likelihood
and the Prior of a
Single Neuron

Mathematical Details: Single Neuron



Learning as Optimization

$$l_r(w) = e(w) + \alpha r(w)$$

$$l_r(w) = - \sum_{i=1}^n \left[y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right] + \frac{\alpha}{2} w^T w$$

From Unregularized Loss to Likelihood

$$l(\omega) = -\sum_{i=1}^n \left[y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right]$$

$$e^{-l(\omega)} = \prod_{i=1}^n \hat{y}^{(i)} \cdot (1-\hat{y}^{(i)})^{(1-y^{(i)})}$$

Remember?

$$\text{Ber}(y|\theta) = \prod_{i=1}^n \theta^{y^{(i)}} \cdot (1-\theta)^{(1-y^{(i)})}$$

$$\underbrace{P(y|x, \omega)}_{\text{Likelihood}} = \prod_i^n \text{Ber}(y^{(i)}|\hat{y}^{(i)}) \rightarrow \text{Model of the NN}$$

$$= \prod_i^n \text{Ber}(y^{(i)}|g(x^{(i)}\omega^\top)) = \prod_i^n g(x^{(i)}\omega^\top)^{y^{(i)}} \cdot (1-g(x^{(i)}\omega^\top))^{(1-y^{(i)})}$$

From Regularizer to Prior

$$\alpha r(\omega) = \frac{\alpha}{2} \omega \omega^T$$

$$e^{-\alpha r(\omega)} = e^{-\frac{\alpha}{2} \omega \omega^T}$$

Remember?

$$N(\omega | \mu=0, \Sigma=\alpha^{-1}) = \frac{1}{Z} e^{-\frac{\alpha}{2} \omega \cdot \omega^T} = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} (\omega - \mu)(\omega - \mu)^T}$$

$$P(\theta) = N(\omega | \mu=0, \Sigma=\alpha^{-1})$$

Gaussian Prior

Posterior over the Weights

$$P(\underline{w} | \underline{x}, \underline{y}) = \frac{\prod_{i=1}^n \text{Ber}(y_i | g(x_i^\top \underline{w})) N(\mu_w, \Sigma_w)}{P(\underline{y})}$$

Bayesian Prediction

$$P(w|x,y) = \frac{\prod_{i=1}^n \text{Ber}(w_i | \sigma(x_i^\top w)) N(\mu_w, \Sigma_w)}{P(y)}$$

$$P(y=1|x^*, D) = \int P(y=1|x^*, w) \cdot P(w|D) dw$$

↑
new input

$$P(y=1|x^*, w) = \sigma(x^* w^\top)$$

$$P(y=0|x^*, w) = 1 - \sigma(x^* w^\top)$$

So far so good...

What's the problem now?

Nasty Likelihood Dist. and Intractable Integrals

$$P(w|x,y) = \frac{\prod_{i=1}^n P(y^{(i)}|x_i, w) \cdot P(w)}{P(y)}$$

Problem 1

Nasty Likelihood

$$\int \prod_{i=1}^n P(y^{(i)}|x_i, w) \cdot P(w) dw$$

$$P(y^{*}=1 | \vec{x}, D) = \int P(y^{*}=1 | \vec{x}, w) \cdot P(w | D) dw$$

Problem 2

Intractable Integrals

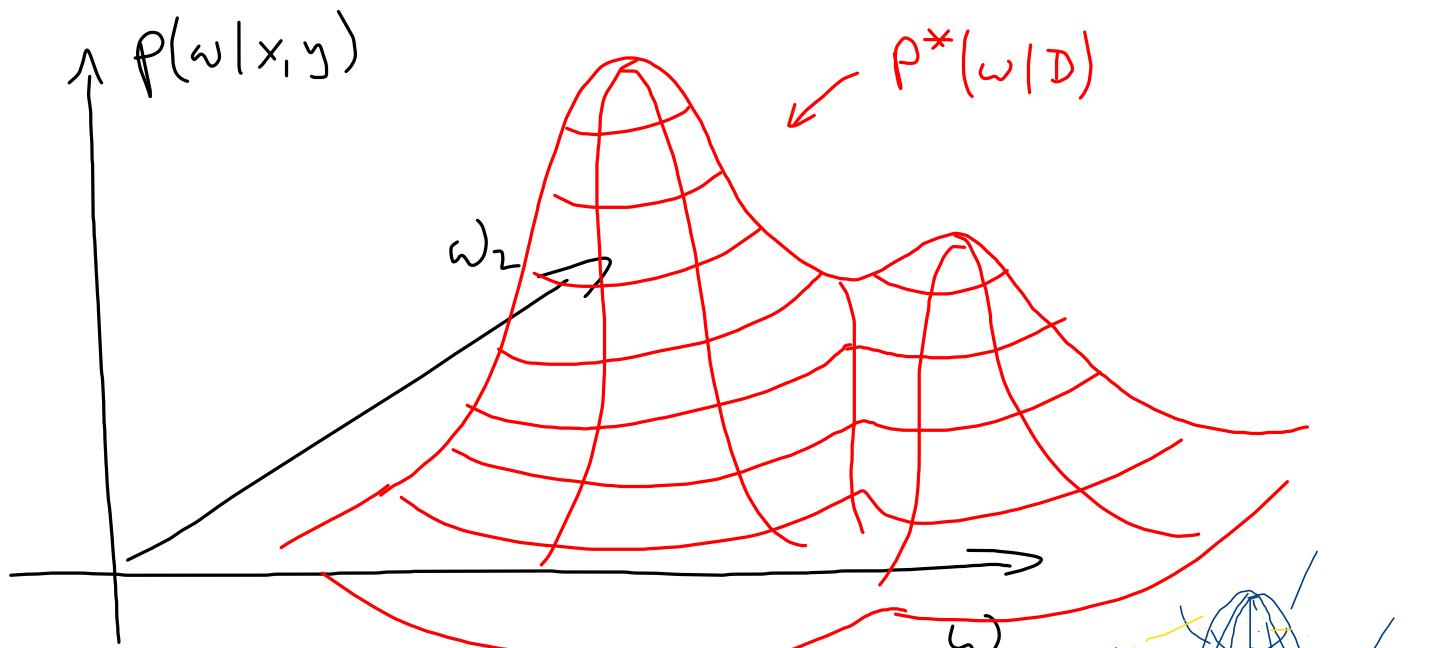
How to solve
an intractable
integral?

From Integral to Sample Mean

$$P(y^*=1 | D) = \int P(y=1 | x^*, w) \cdot P(w | D) dw$$

$$= E \left[y(x^*, \omega) \right]_{\omega \sim P(\omega|D)}$$

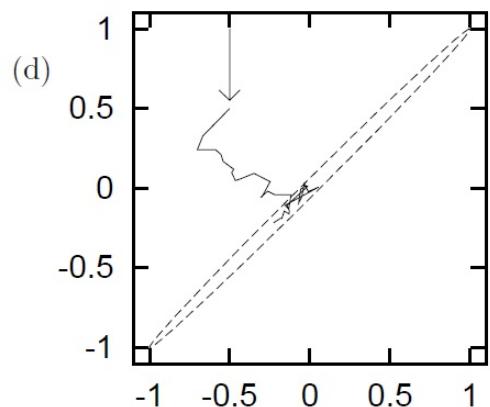
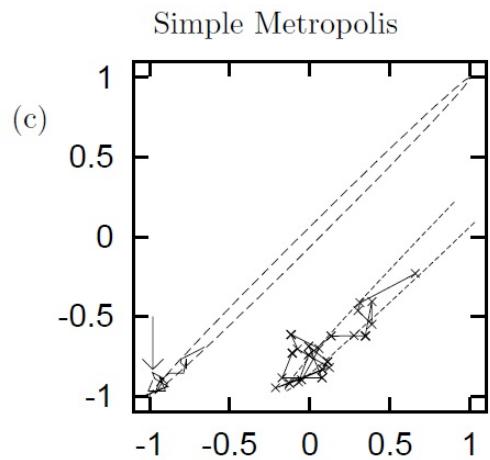
Sampling from the Posterior: MCMC



$$P(\omega | D) = \frac{P^*(\omega | D)}{Z}$$



Random Walk Visualization



Python Code: Weight Samples

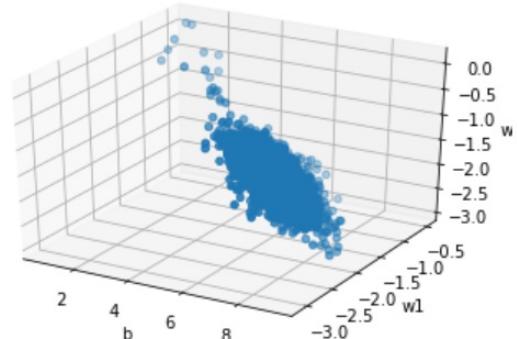
```
def MCMC(function, X, Y, N, d, sigma=0.1):
    mu0 = [0]*d
    cov = np.diag([sigma]*d)
    w_old = np.random.multivariate_normal(mu0, cov)
    sample_list = []
    for i in range(N):
        w_new = np.random.multivariate_normal(w_old, cov)
        a = np.exp(function(X, Y, w_new) - function(X, Y, w_old))
        u = np.random.uniform(0, 1)
        if (a > u):
            sample_list.append(w_new)
        w_old = w_new
    return np.array(sample_list)
```

```
def sigmoid_layer(X, W):
    x_biases = np.full((X.shape[0], 1), 1)
    X_ = np.concatenate((x_biases, X), axis=1)
    odds = np.exp(-np.dot(X_, W.T))
    y_hat = 1/(1 + odds)
    return y_hat
```

```
def log_posterior_over_parameters_logreg(X, Y, W):
    sigma_prior = 1.5
    log_prior = (-np.dot(W, W.T)/(2*sigma_prior**2))
    y_hat = sigmoid_layer(X, W)
    full_log_likelihood = (Y.T)*np.log(y_hat) + (1 - Y.T)*np.log(1-y_hat)
    post = log_prior + np.sum(full_log_likelihood)
    return post
```

```
W_samples = (MCMC(log_posterior_over_parameters_logreg, X, Y, N=10000, d=3))
```

$$\begin{aligned} &\cancel{\log p(w_{new}|D)} \\ &\cancel{\log p(w_{old}|D)} \end{aligned}$$



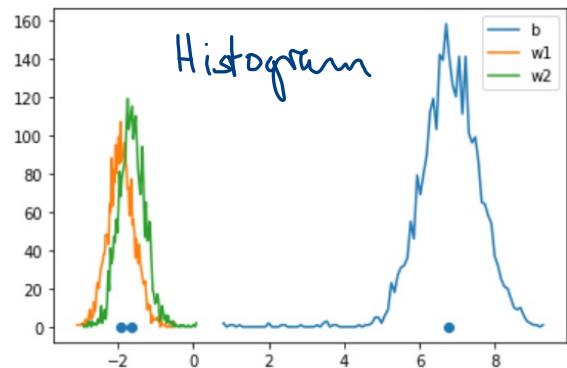
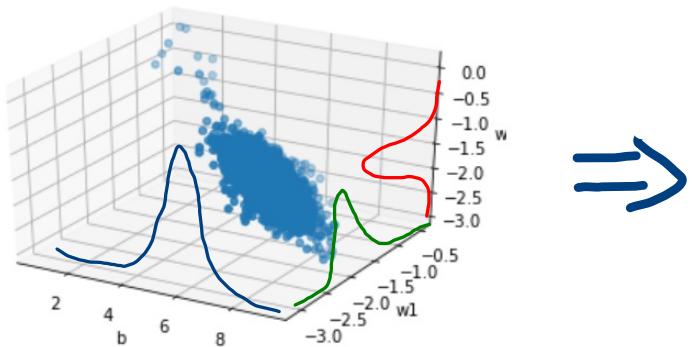
$$\hat{y} = \frac{1}{1 + e^{-X\omega^T}}$$

↙ odds

$$\log(e^{-\frac{1}{2}\sigma^2\omega\omega^T})$$

$$\log^*(y|D) = \sum \log[\text{Ber}(y|\hat{y})] + \log[N(\omega|\mu, \sigma^2)] = \text{Loss function}$$

Visualization as Histogram



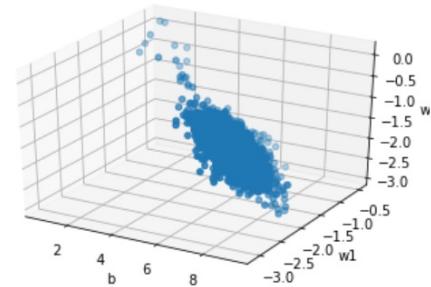
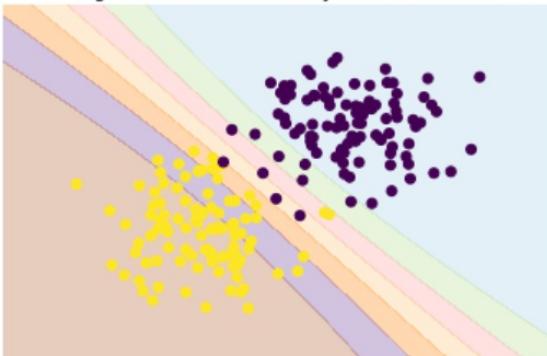
Python Code : Prediction

$$P(y^*|D) = \frac{1}{n} \sum_{i=1}^n \tilde{G}(x_i^*, w^{(i)})$$

```
def predict_MC(X, W_samples):
    return (np.nanmean(sigmoid_layer(X, W_samples), axis=1)/W_samples.shape[0])
```

```
y_pred_full = predict_MC(X_pred, W_samples)
```

Single Neuron with Full Bayesian Prediction



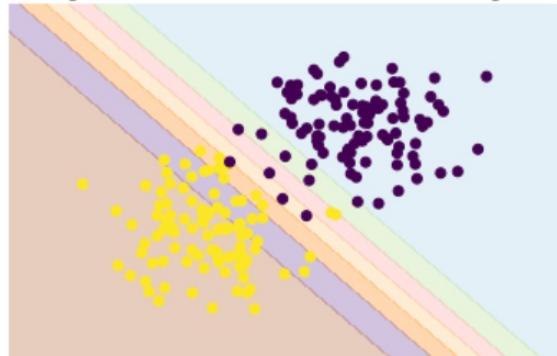
$$\bar{w}_{MAP} = \frac{1}{N} \sum_{i=1}^N w^{(i)}$$

$$\hat{y} = \tilde{G}(x \bar{w}_{MAP}^T)$$

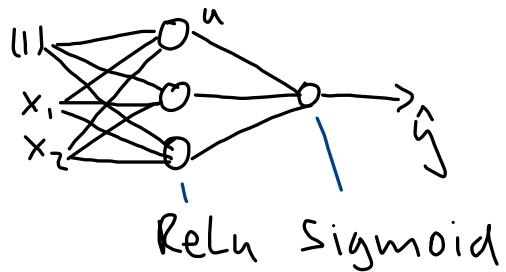
POINT ESTIMATE

```
W_map = np.median(W_samples, axis=0).reshape(1, -1)
y_pred_map = sigmoid_layer(X_pred, W_map)
```

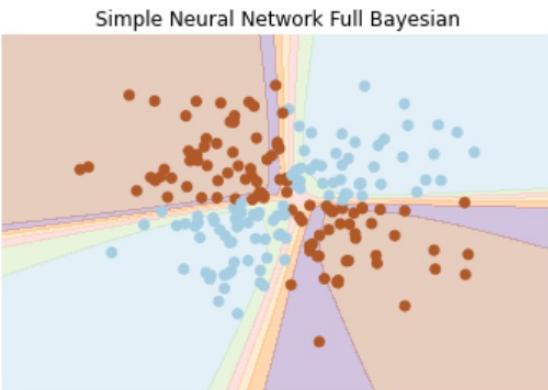
Single Neuron with Point Estimates of the Weights



XOR Data with Simple Neural Network



$$\hat{y} = \sigma(\max(0, x \cdot w_1^\top) \cdot w_2^\top)$$



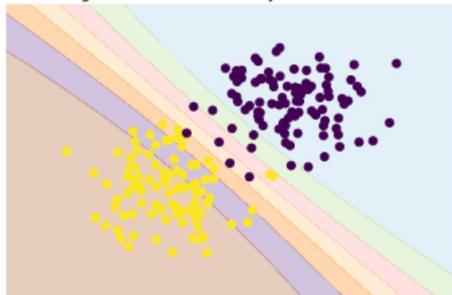
Simple Neural Network with Point Estimates of Weights



Dependence of Uncertainty on Sample Size

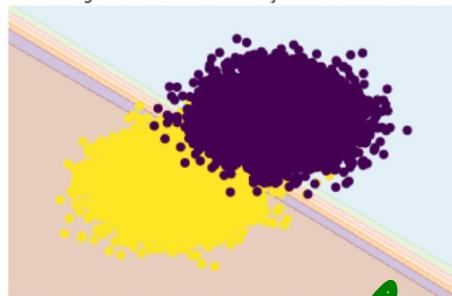
$n=200$

Single Neuron with Full Bayesian Prediction

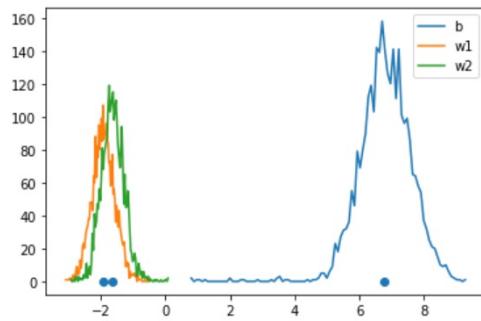


$n=20000$

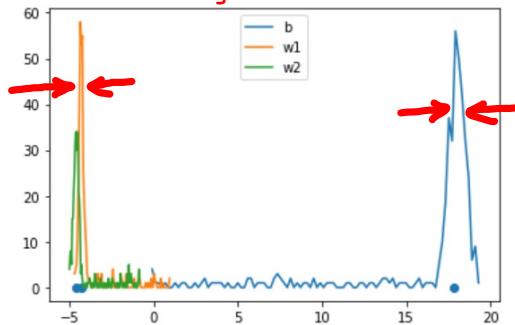
Single Neuron with Full Bayesian Prediction



Looks like
 ω_{MLE}

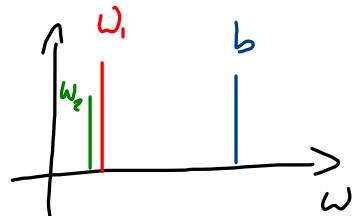


Uncertainty shrinks if $n \rightarrow \infty$



$$\lim_{n \rightarrow \infty} \omega_{MLE} = \omega_{MAP}$$

$$\lim_{n \rightarrow \infty} G_w = 0$$



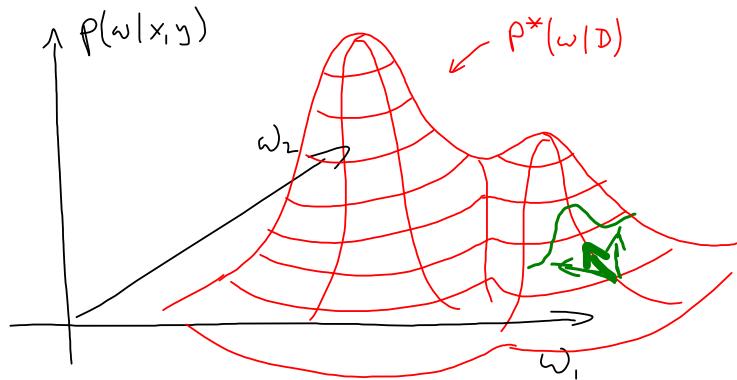
Disadvantages of Metropolis Algorithm

- Many rejected samples
 - Can wander off in many directions
- } computationally expensive
→ slow

Future Work

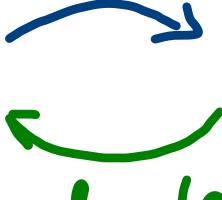
→ Hamiltonian Monte Carlo
("Hybrid")

Optimization and Random Walk Hybrid



→ Speed

Take Away

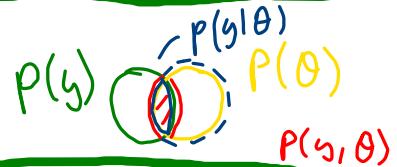
- Black Swan Paradox
 $e^{-\text{loss}}$
- Loss  Probability
 $-\log(p)$
- $\int f(x) \cdot p(x) dx = E[f(x)] \underset{x \sim p(x)}{\approx} \frac{1}{n} \sum_i^n f(x^{(i)})$
- $\lim_{n \rightarrow \infty} \omega_{MLE} = \omega_{MAP} \Rightarrow \lim_{n \rightarrow \infty} \zeta_w = 0$

Understanding
Numerator and Denominator
of Bayes

(Optional?)

Rules of Probability

$P(\mathcal{S})$



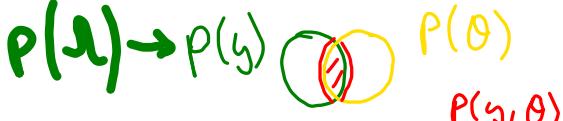
Creates a joint probability

Product Rule

$$P(y, \theta) = p(y|\theta) P(\theta)$$

$$p(\theta|y) = \frac{p(y|\theta) P(\theta)}{p(y)} = \frac{P(y, \theta)}{p(y)}$$

All possible y 's



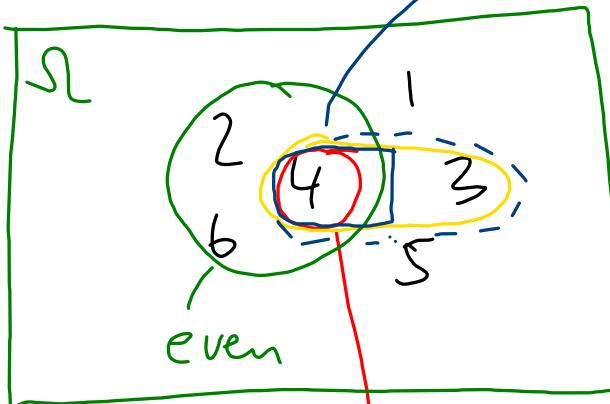
The reference system switches from $\mathcal{S} \rightarrow y$

Example



Die

$$P(\text{even} | d12) = \frac{1}{2}$$



$$P(\text{even}) = \frac{1}{2} \quad \text{both} \quad P(\text{even}, d12) = \frac{1}{6} \leftarrow \begin{matrix} \#4 \\ |\Omega| \end{matrix}$$

$$P(d12 | \text{even}) = \frac{P(\text{even}, d12)}{P(\text{even})} = \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2}} = \frac{1}{3} //$$

$$\overbrace{P(\text{even} | d12) \cdot P(d12) + P(\text{even} | \text{not } d12) \cdot P(\text{not } d12)}$$