

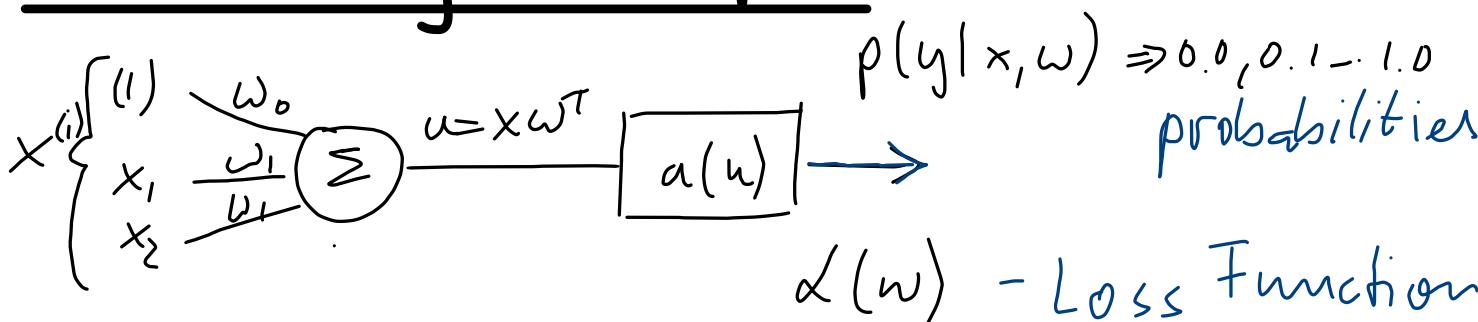
Bayesian Machine Learning
&

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

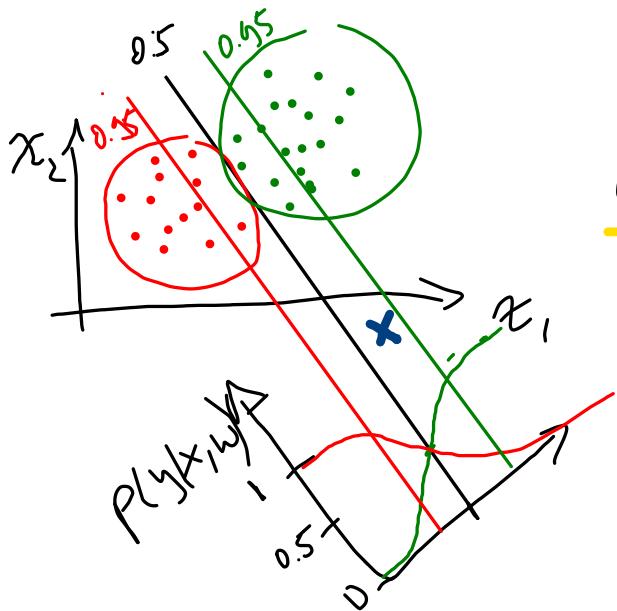


$\ell(\omega)$ - Loss Function

$\nabla_{\omega} \ell(\omega)$ - Gradient

$$\underline{\omega} \leftarrow \omega - \alpha \nabla_{\omega} \ell(\omega)$$

Gradient descent



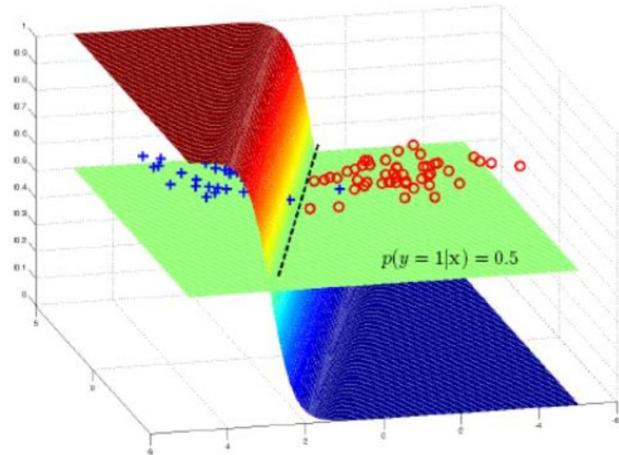
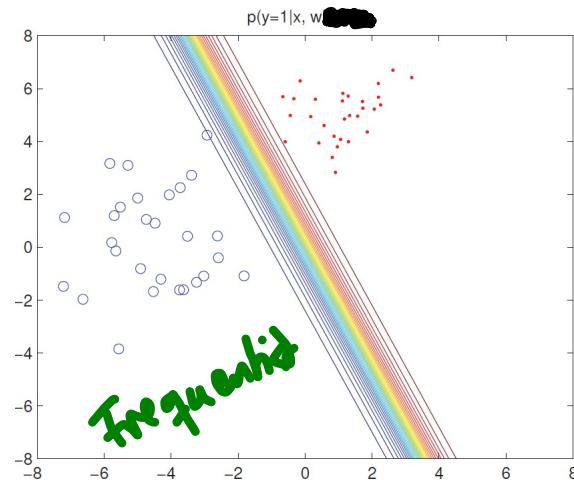
Frequentist
Point Estimation

Motivating Example

$$x^{(i)} \left\{ \begin{array}{l} (1) \\ x_1 \\ x_2 \end{array} \right. \xrightarrow{\omega_0} \sum \xrightarrow{u = x\omega^T} \boxed{a(u)} \xrightarrow{} p(y|x, \omega) \Rightarrow 0.0, 0.1 - 1.0 \text{ probabilities}$$

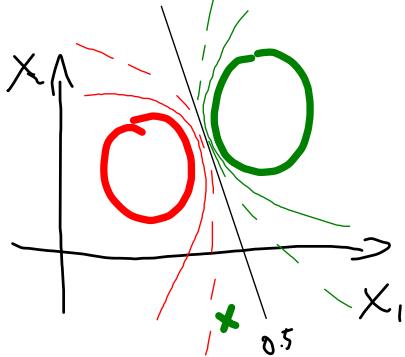
$\alpha(w)$ - Loss Function

$\nabla_w \alpha(w)$ - Gradient

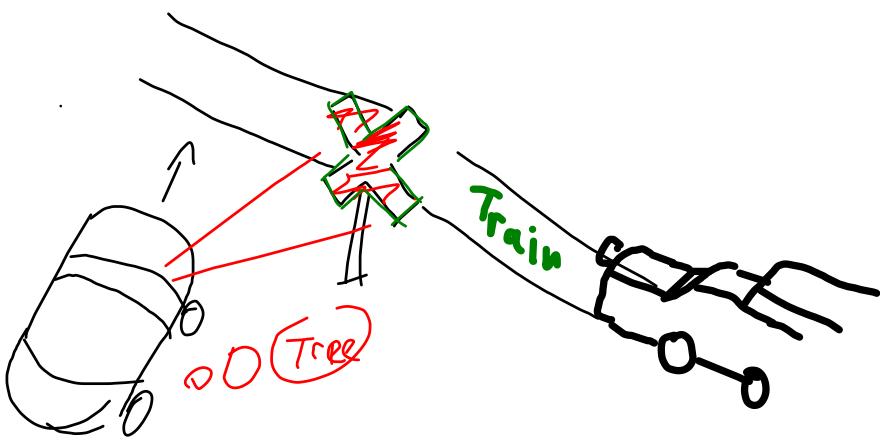
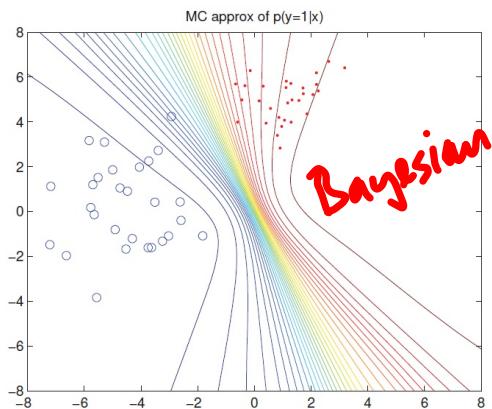


$\alpha(w)$
t descent)

Uncertainty: Single Neuron \rightarrow Bayesian

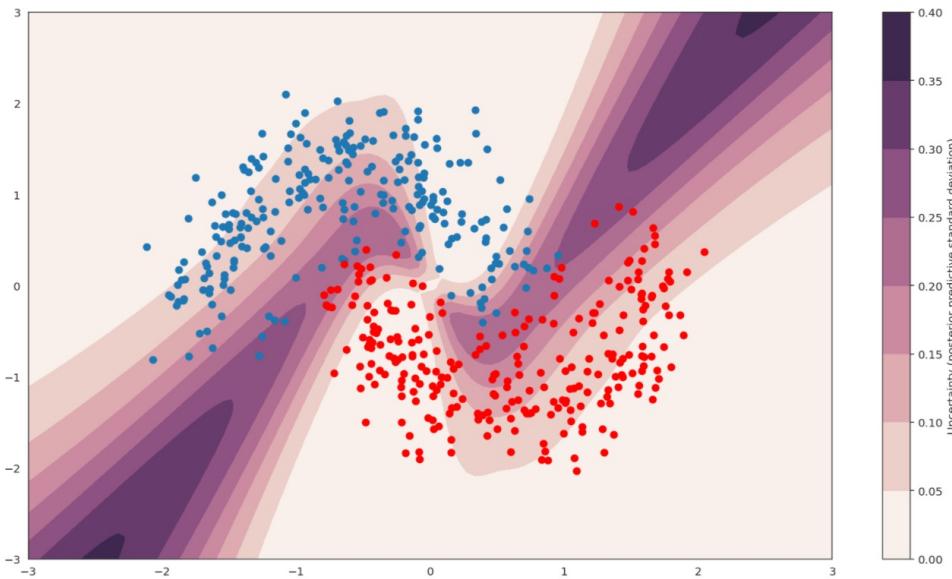
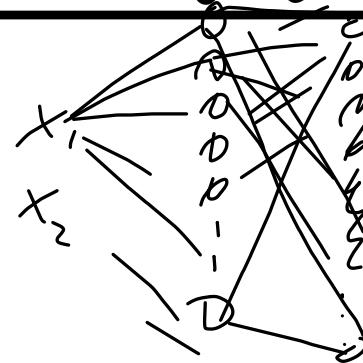


- Medical field
- Automotive

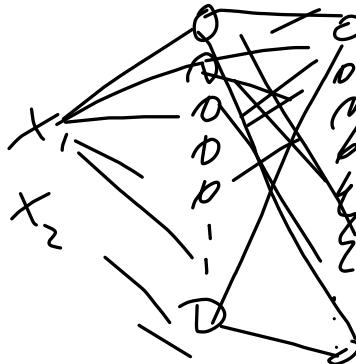
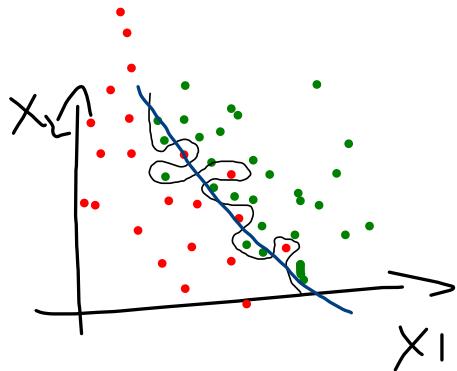


More Complex Data and Model

Neural Net



Over fitting



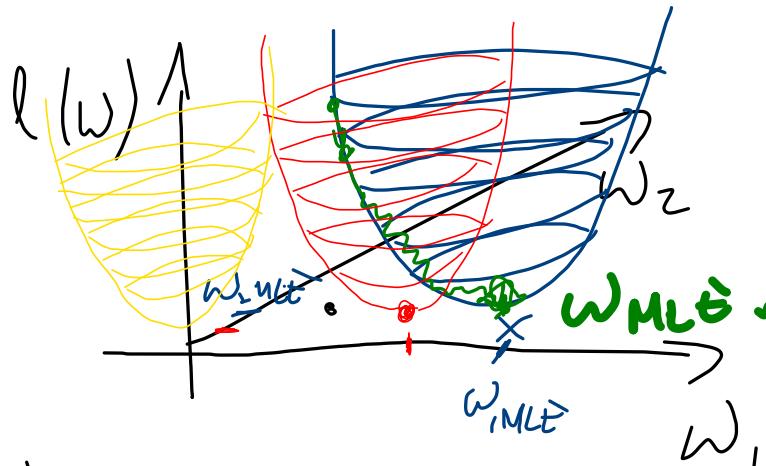
Optimization World

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

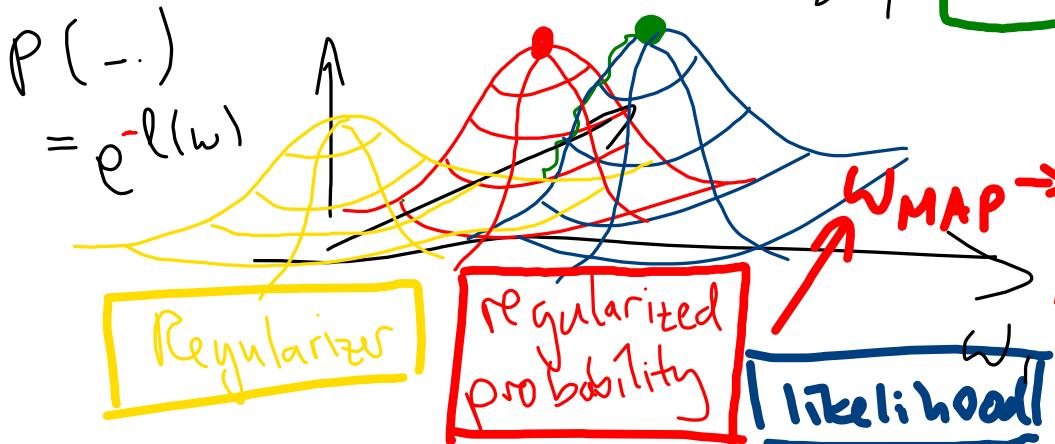
Regularizer

Bayesian Point Estimate

Point Estimation vs. Probability Distribution over Parameters (Weights)



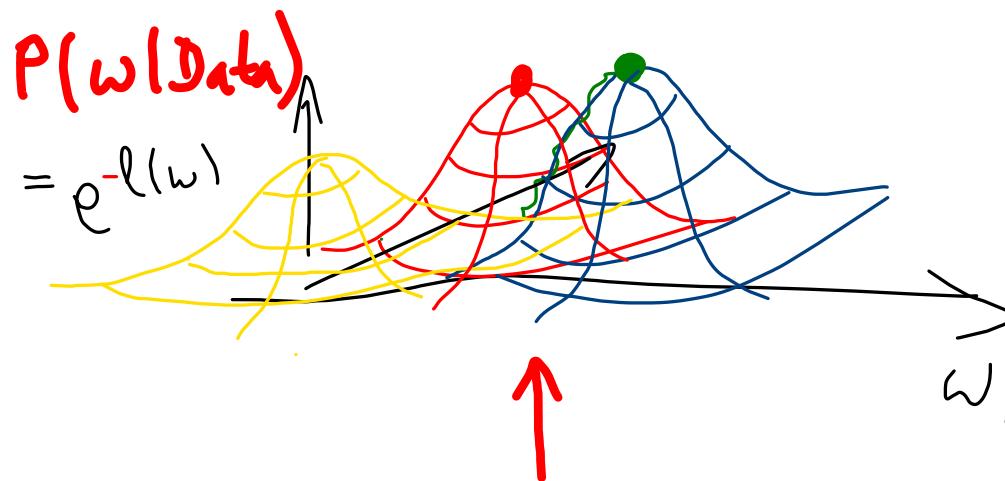
$w_{MLE} \rightsquigarrow$ Frequentist Point Estimate



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

Regularized Logistic Regression

Goal: Probability Distribution over Parameters

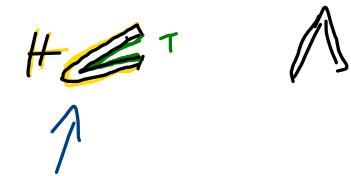


We want the whole distribution
over the weight,

instead of a point estimate

Bayesian Probability Intuition

Biased Coin Flips with „Small Data“



bend
coin

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

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“

H ↗ A ↘

$y = \{1, 0\}$

„Hallucinations“
probability of getting Heads

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

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

$$\underset{i=1}{\overset{n}{\times}} \text{Ber}(y_i | \theta)$$

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

"?"

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

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

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

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

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

} Conjugate Prior

Prior - Posterior

deterministic

$$\theta_i$$



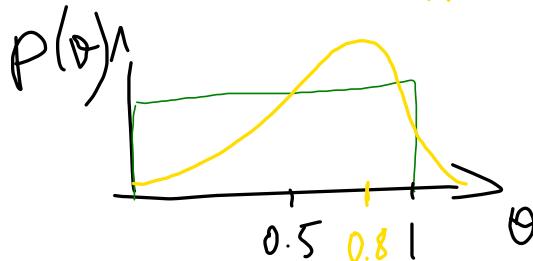
$$\boxed{\arg \max_{\theta} \text{Ber}(y|\theta)}$$

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

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

random
prior:

i.e. Beta(α, β)



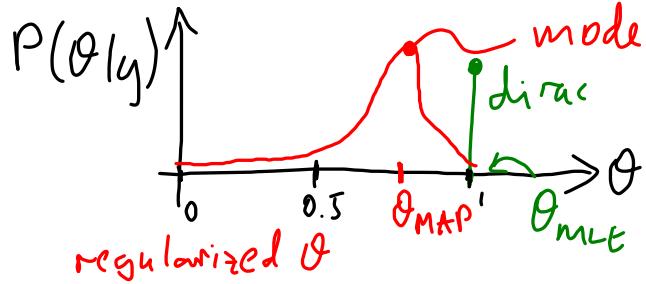
Like likelihood

Norm.

Posterior

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

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



General Bayes Formula

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

(Model) likelihood prior (Hallucinations)

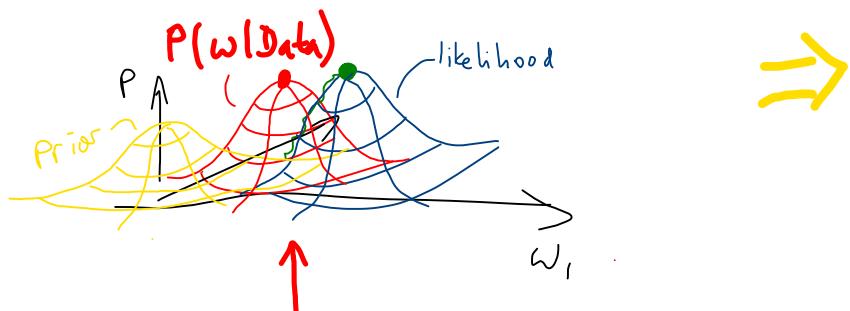
Posterior over parameters (weights)

normalization constant
aka. prob. of data

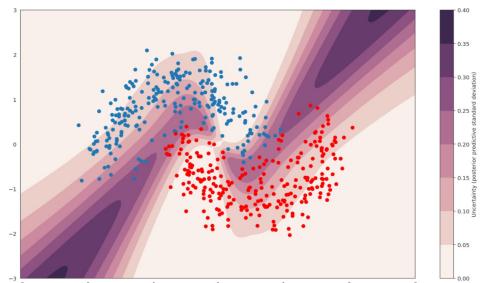
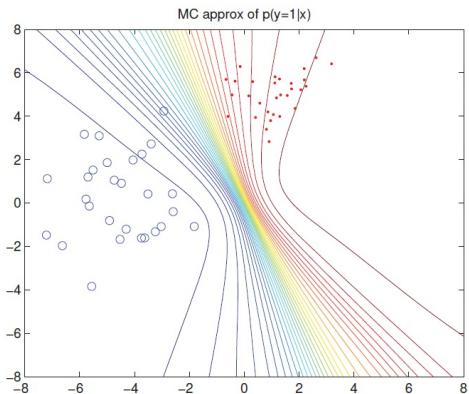
→ makes it sum to one

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

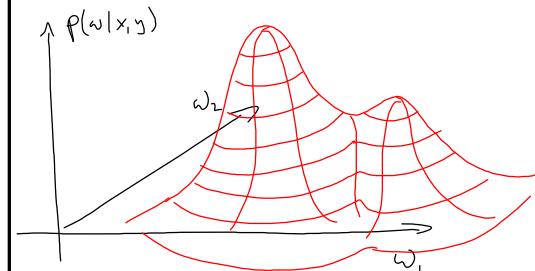
Multiclass Data • Dirichlet-Categorical

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

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

?

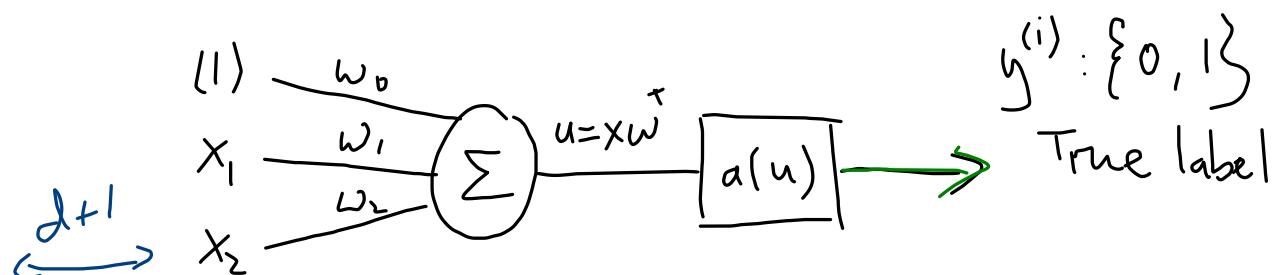
prior (Hallucinations)

↑

Also adding features x

↖ normalization const.
integral of numerator

Recap: Single Neuron Optimization



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

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

Learning as Optimization

Regularized
Loss Function

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

Gradient
Update

$$\underbrace{\tilde{V}_w \ell(w)}_{w \leftarrow w - \alpha \tilde{V}_w \ell(w)}$$

From Regularized Loss Function to Bayes

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

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

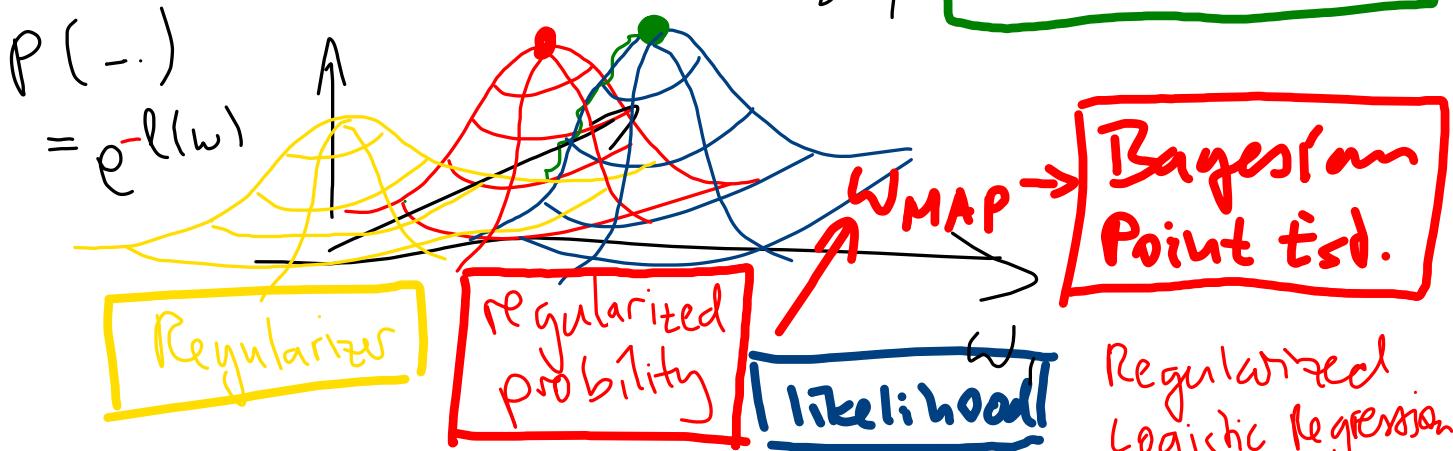
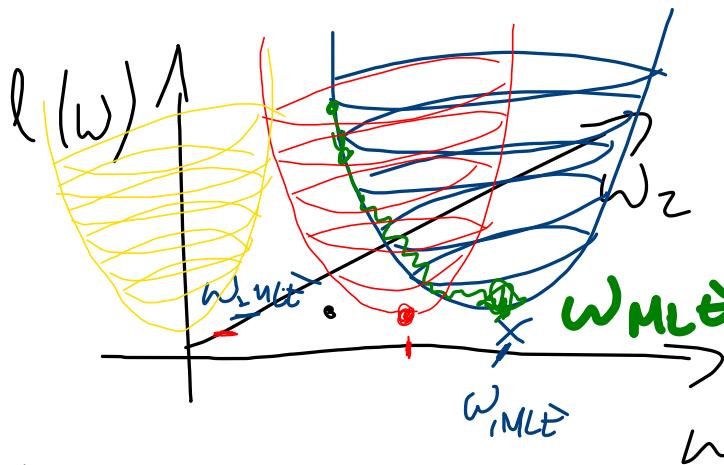
(Model) likelihood prior (Hallucinations)

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

Z ↗ normalization const.
integral of numerator

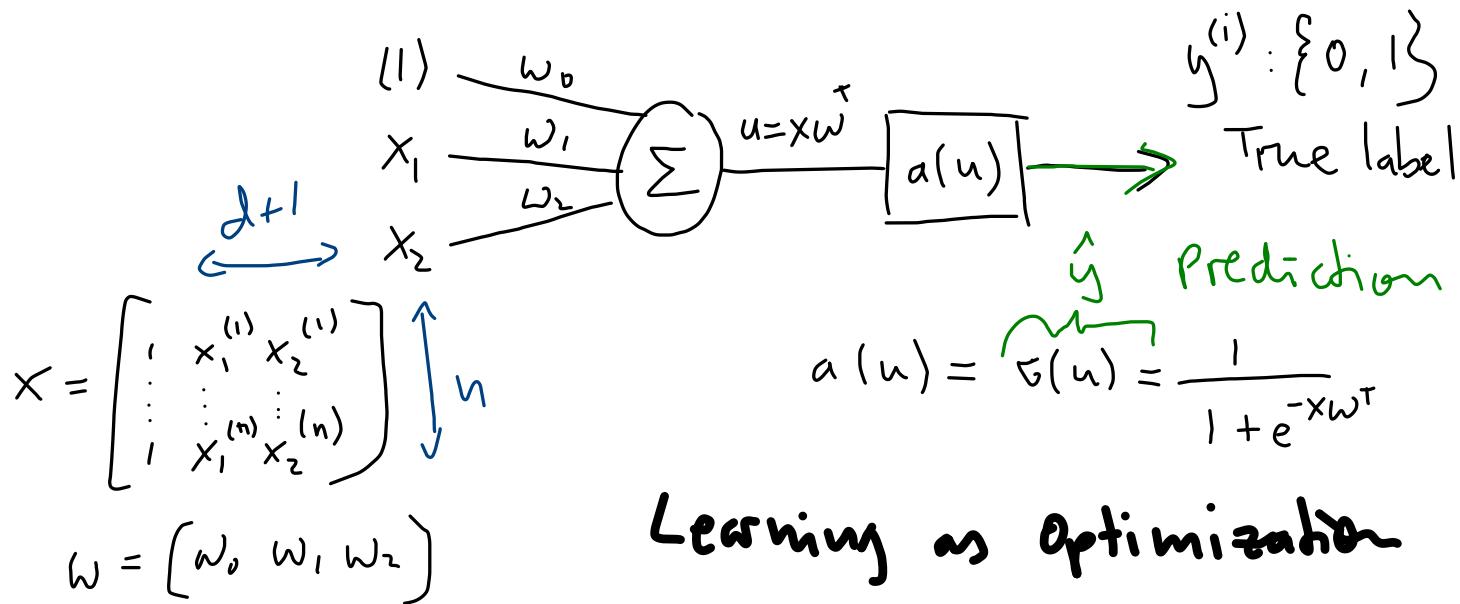
$$\int_{-\infty}^{\infty} e^{-l(\omega)} \cdot e^{-\alpha r(\omega)} d\omega$$

Remember this Connection?



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?

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

$$\underbrace{P(y|x, \omega)}_{\text{Likelihood}} = \text{Ber}(y^{(i)}|\hat{y}^{(i)})$$
$$= \text{Ber}(y|g(x\omega^\top))$$

From Regularizer to Prior

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

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

Remember?

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

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

Gaussian Prior

Posterior over the Weights

$$P(\omega | x, y) = \frac{\text{Ber}(y | g(x \omega^\top)) N(\mu_\omega, \Sigma_\omega)}{P(y)}$$

Bayesian Prediction

$$P(w|x,y) = \frac{\text{Ber}(w|\mathcal{G}(x|w^\top)) N(\mu_w, \sigma_w)}{P(y)}$$

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

↑
new input

So far so good...

What's the problem now?

Nasty Likelihood Dist. and Intractable Integrals

$$P(w|x,y) = \frac{P(y|x,w) \cdot P(w)}{P(y)}$$

Problem 1

Nasty Likelihood

$$\int_{-\infty}^{\infty} P(y|x,w) \cdot P(w) dw$$

Problem 2

Intractable Integrals

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

How to solve
an intractable
integral?

From Integral to Sample Mean

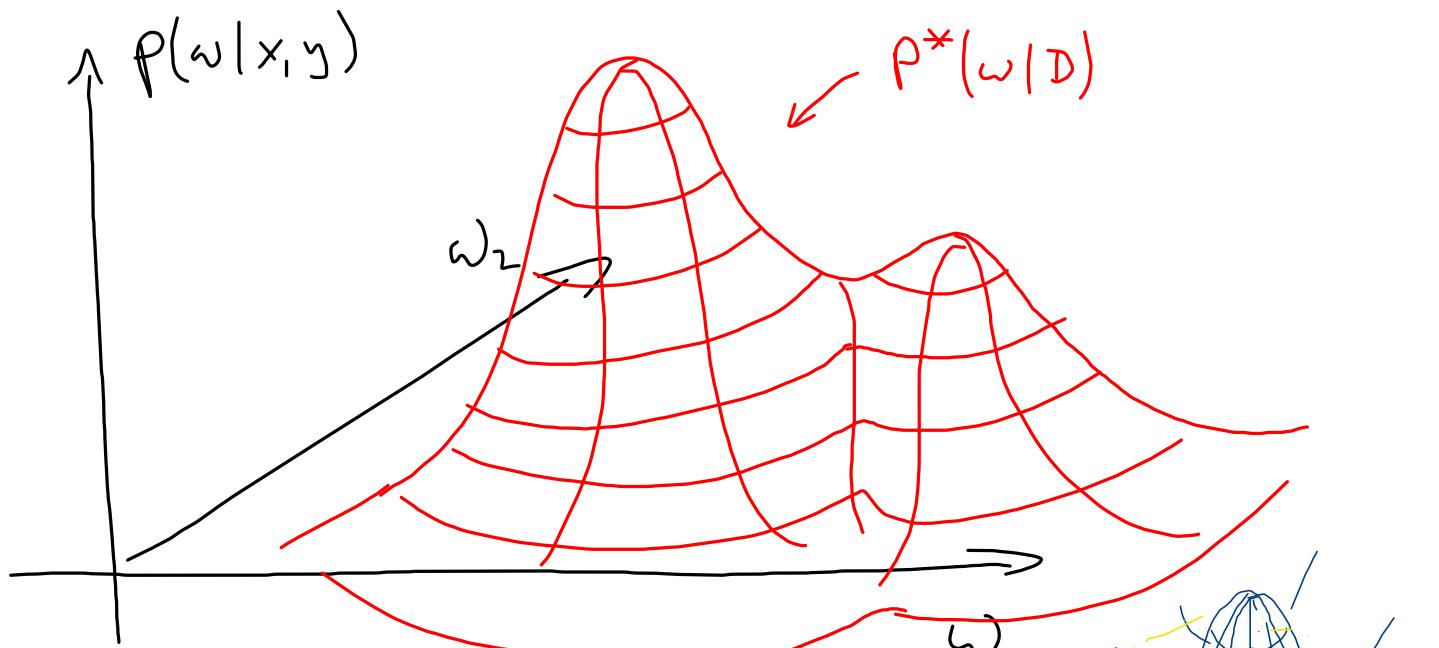
$$P(y^*|D) = \int P(y|x^*, w) \cdot P(w|D) dw$$

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

$$\approx \frac{1}{n} \sum_{i=1}^n y(x^*, \omega_i) \quad \text{where } \omega_i \sim p(\omega | D)$$

↑ ↑
 given needed

Sampling from the Posterior: MCMC



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



Pseudocode of MC MC

Algorithm 24.2: Metropolis Hastings algorithm

1 Initialize x^0 ; $\leftarrow \omega$

2 **for** $s = 0, 1, 2, \dots$ **do**

3 Define $x = x^s$;

4 Sample $x' \sim q(x'|x)$;

5 Compute acceptance probability

$$\alpha = \frac{\tilde{p}(x')q(x|x')}{\tilde{p}(x)q(x'|x)} \text{ if Gaussian}$$

6 Compute $r = \min(1, \alpha)$;

7 Sample $u \sim U(0, 1)$;

8 Set new sample to

$$x^{s+1} = \begin{cases} x' & \text{if } u < r \\ x^s & \text{if } u \geq r \end{cases}$$

here $X = \omega$

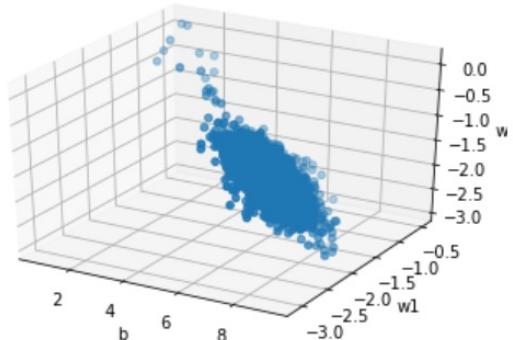
Python Code

```
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)
    odds1 = np.exp(np.dot(X_, W.T))
    U = odds1/(1 + odds1)
    return U

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

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

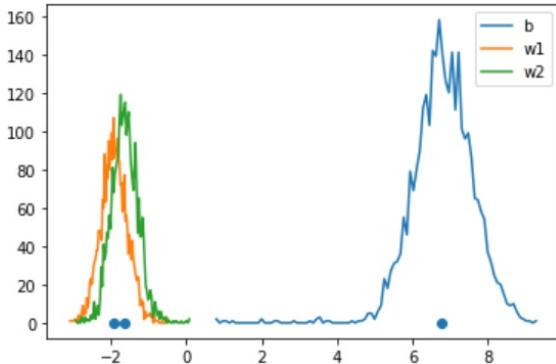
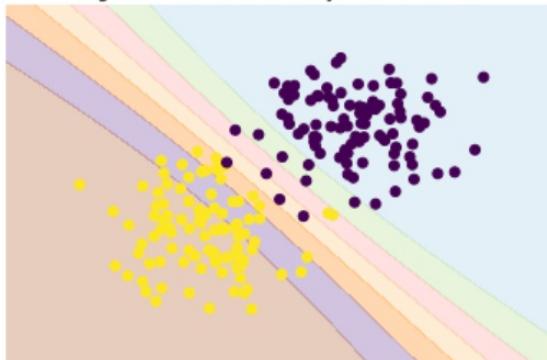


Python Code

```
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)
```

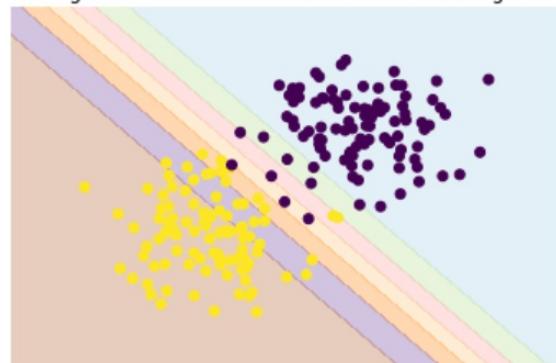
Single Neuron with Full Bayesian Prediction



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

Simple Neural Network Full Bayesian



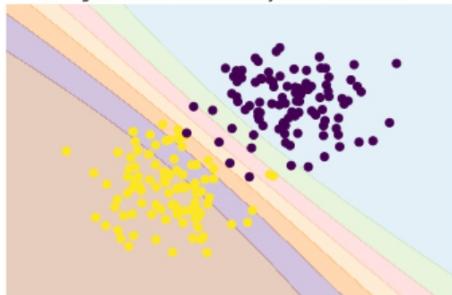
Simple Neural Network with Point Estimates of Weights



Dependence of Uncertainty on Sample Size

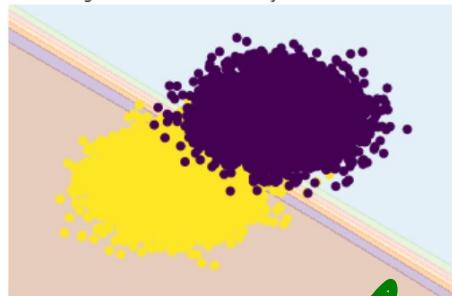
$n=200$

Single Neuron with Full Bayesian Prediction

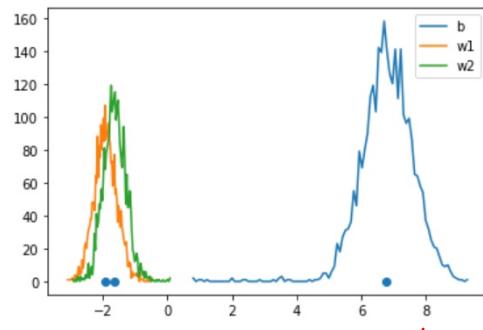


$n=20000$

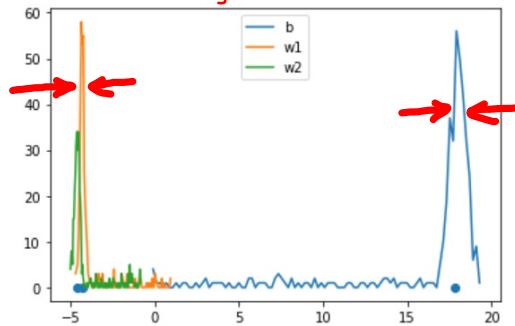
Single Neuron with Full Bayesian Prediction



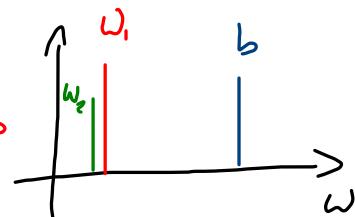
Looks like
 ω_{MLE}



Uncertainty shrinks if $n \rightarrow \infty$



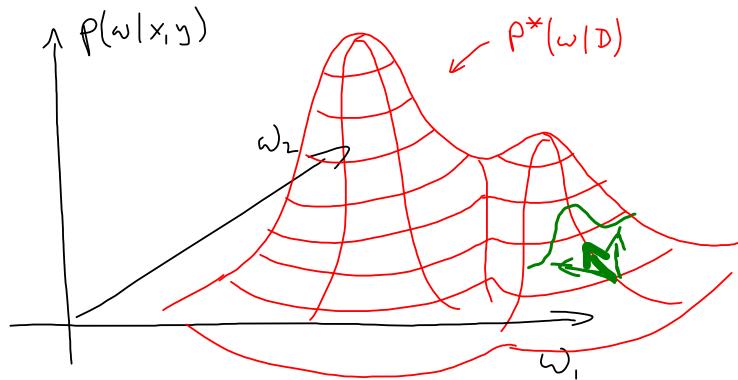
$$\lim_{n \rightarrow \infty} \omega_{MLE} = \omega_{MAP}$$



Future Work

→ Hamiltonian Monte Carlo
("Hybrid")

Optimization and Random Walk Hybrid



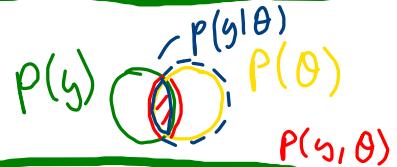
→ Speed

Understanding
Numerator and Denominator
of Bayes

(Optional?)

Rules of Probability

$P(\mathcal{S})$



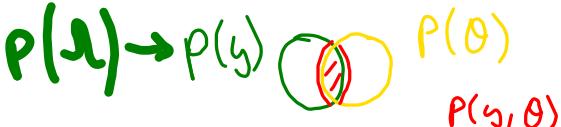
Creates a joint probability

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



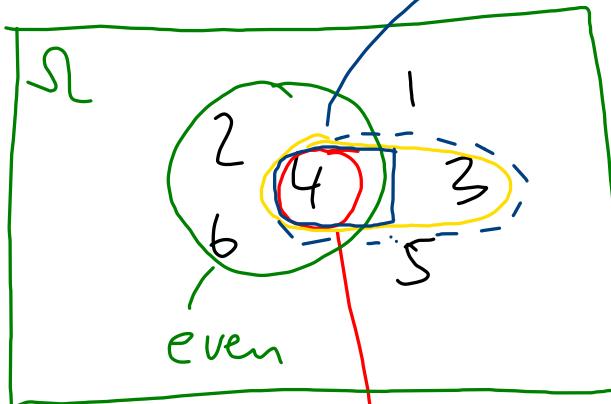
The reference system switches from $\mathcal{S} \rightarrow y$

Example



Die

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



dividers of 12
"d12"

$$P(d12) = \frac{1}{3}$$

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

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

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