


ISTA 421 + INFO 521
**Introduction to
Machine Learning**

Lecture 9:
**Maximum Likelihood –
Uncertainty in Parameters
and Predictions**

Clay Morrison
claytonm@email.arizona.edu

20 September 2017

 1

Topics

- Maximum Likelihood
- (What to expect from) Expectation
- The generative picture
- Maximum Likelihood Estimation
 - Uncertainty in parameters
 - Uncertainty in predictions

Explicitly Modeling Uncertainty

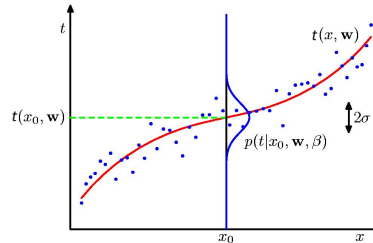
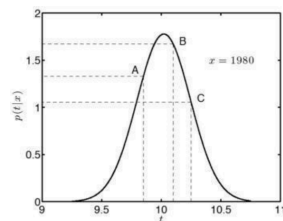
$t_n = f(\mathbf{x}_n; \mathbf{w}) + \epsilon_n, \epsilon_n \sim \mathcal{N}(0, \sigma^2)$ Adding a R.V to our linear model: ... to model noise/variance/uncertainty in the response (t).

$p(t_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$ Now our "response" is itself a random variable. (Here we assume Gaussian/Normal noise.)

We assume (conditional) **independence**: given a particular input \mathbf{x} and parameters \mathbf{w} , the predicted variance at that point is indep. of any other point. However, we also assume the variance form (i.e., what type of distribution) and its parameters will be **identical** at any point (so are constant, not a fn of \mathbf{x} or \mathbf{w}). --Hence, i.i.d.

$$L = p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$$

We can now use the data (input and target pairs) to give us a "likelihood" score of how well the model "fits" the mass/density of our model's uncertainty over the data. The more the data all fits under the bulk of the probability mass/density, the more it looks like our model could have generated the data (the more "likely" the data looks given our model).



Maximize the Likelihood

$$L = p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$$

Since we are working with a product of Gaussians, which in turn include The exponential function (e), take the natural log (often just represented generically as $\log(L)$)

$$\begin{aligned} \log L &= \sum_{n=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma^2} \exp \left\{ -\frac{1}{2\sigma^2} (t_n - f(\mathbf{x}_n; \mathbf{w}))^2 \right\} \right) \\ &= \sum_{n=1}^N \left(-\frac{1}{2} \log(2\pi) - \log \sigma - \frac{1}{2\sigma^2} (t_n - f(\mathbf{x}_n; \mathbf{w}))^2 \right) \\ &= -\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - f(\mathbf{x}_n; \mathbf{w}))^2. \end{aligned}$$



$$\frac{df}{dx} = \frac{df}{dg} \cdot \frac{dg}{dx}$$

Maximize the Likelihood: \mathbf{w}

$$\begin{aligned}\log L &= -\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - f(\mathbf{x}_n; \mathbf{w}))^2 \\ &= -\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2\end{aligned}$$

$$\begin{aligned}\frac{\partial \log L}{\partial \mathbf{w}} &= \frac{1}{\sigma^2} \sum_{n=1}^N \mathbf{x}_n (t_n - \mathbf{x}_n^\top \mathbf{w}) \\ &= \frac{1}{\sigma^2} \sum_{n=1}^N \mathbf{x}_n t_n - \mathbf{x}_n \mathbf{x}_n^\top \mathbf{w} = \mathbf{0}\end{aligned}$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}, \quad \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}$$

$$\begin{aligned}\frac{\partial \log L}{\partial \mathbf{w}} &= \frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{t} - \mathbf{X}^\top \mathbf{X} \mathbf{w}) = \mathbf{0} & \frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{t} - \mathbf{X}^\top \mathbf{X} \mathbf{w}) &= \mathbf{0} \\ & & \mathbf{X}^\top \mathbf{t} - \mathbf{X}^\top \mathbf{X} \mathbf{w} &= \mathbf{0} \\ & & \mathbf{X}^\top \mathbf{X} \mathbf{w} &= \mathbf{X}^\top \mathbf{t} \\ & & \mathbf{w} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}\end{aligned}$$

Maximize the Likelihood: σ

$$\log L = -\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

$$\frac{\partial L}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{n=1}^N (t_n - \mathbf{x}_n^\top \hat{\mathbf{w}})^2 = 0$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{x}_n^\top \hat{\mathbf{w}})^2$$

$$\begin{aligned}\sigma^2 &= \frac{1}{N} (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})^\top (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}}) \\ &= \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - 2\mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}} + \hat{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{X} \hat{\mathbf{w}})\end{aligned}$$

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - 2\mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} + \mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}) \\ &= \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - 2\mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} + \mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}) \\ &= \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t})\end{aligned}$$

$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

Simplify further by plugging in

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

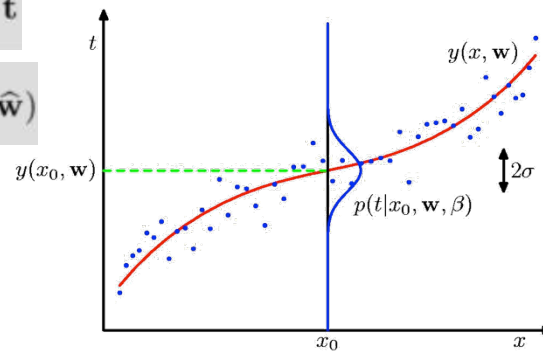
Maximum Likelihood

$$L = p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

Predictive distribution



The equations are unique

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

Hessian Matrix

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_K} \\ \frac{\partial^2 f}{\partial w_2 \partial w_1} & \frac{\partial^2 f}{\partial w_2^2} & \cdots & \frac{\partial^2 f}{\partial w_2 \partial w_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_K \partial w_1} & \frac{\partial^2 f}{\partial w_K \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_K^2} \end{bmatrix}$$

$$\log L = -\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

$$\frac{\partial \log L}{\partial \mathbf{w}} = \frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{t} - \mathbf{X}^\top \mathbf{X} \mathbf{w})$$

$$\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} = -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

This matrix is definite negative, so maximum likelihood solution (with linear model with additive Gaussian noise) is unique. $\mathbf{x}^\top \mathbf{H} \mathbf{x} < 0$

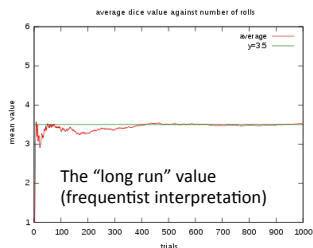


Expectation (a reminder)

- The **expected value** of a random variable is the weighted (by probability) average of all possible values

$$\mathbf{E}_{P(x)} \{f(x)\} = \sum_x f(x)P(x)$$

$$\mathbf{E}_{p(\mathbf{x})} \{f(\mathbf{x})\} = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$



"Belief" in possible worlds
(Bayesian interpretation)

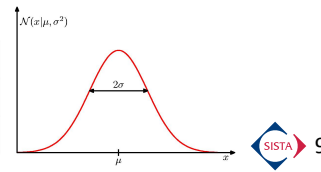
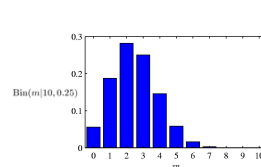
Mean:

$$\mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\} = \sum_{\mathbf{x}} \mathbf{x}P(\mathbf{x})$$

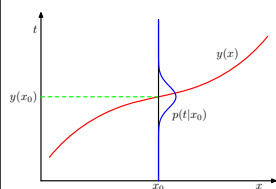
"(Co)Variance":

$$\text{cov} \{\mathbf{x}\} = \mathbf{E}_{P(\mathbf{x})} \{(\mathbf{x} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\})(\mathbf{x} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\})^T\}$$

$$\text{cov} \{\mathbf{x}\} = \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\mathbf{x}^T\} - \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\} \mathbf{E}_{P(\mathbf{x})} \{\mathbf{x}\}^T$$

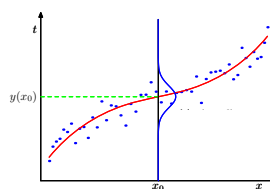


Explicitly Modeling Uncertainty: The Generative Picture

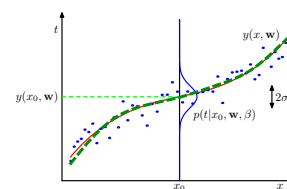


The generating process...

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$



... generates data ...



... that we fit a model to

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \sigma^2)$$

$$= \prod_{n=1}^N \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2)$$

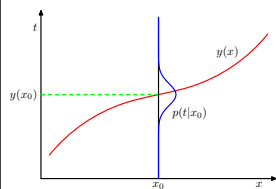
estimated parameters

$$p(\hat{\mathbf{t}}|\mathbf{X}, \hat{\mathbf{w}}, \hat{\sigma}^2) = \prod_{n=1}^N p(\hat{t}_n|\mathbf{x}_n, \hat{\mathbf{w}}, \hat{\sigma}^2)$$

prediction

$$= \prod_{n=1}^N \mathcal{N}(\hat{\mathbf{w}}^T \mathbf{x}_n, \hat{\sigma}^2)$$

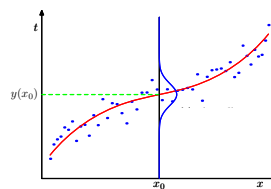
Explicitly Modeling Uncertainty: The Generative Picture



The generating process...

$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\begin{aligned} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) &= \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) \\ &= \prod_{n=1}^N \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2) \end{aligned}$$



... generates data ...

If we generate new data from the generating process, we get different parameters $\hat{\mathbf{w}}$ and $\hat{\sigma}^2$, and a different predictive distribution $p(\hat{\mathbf{t}}|\dots)$

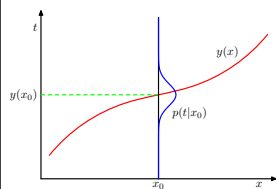
$$\begin{aligned} p(\hat{\mathbf{t}}|\mathbf{X}, \hat{\mathbf{w}}, \hat{\sigma}^2) &= \prod_{n=1}^N p(\hat{t}_n | \mathbf{x}_n, \hat{\mathbf{w}}, \hat{\sigma}^2) \\ &= \prod_{n=1}^N \mathcal{N}(\hat{\mathbf{w}}^\top \mathbf{x}_n, \hat{\sigma}^2) \end{aligned}$$

prediction



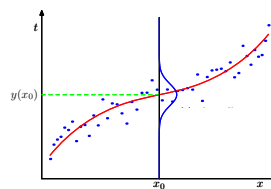
11

Explicitly Modeling Uncertainty: The Generative Picture



The generating process...

$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$



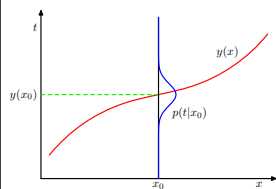
... generates data ...

Suppose we generate 10,000 datasets (each slightly different) but from the same generating process...



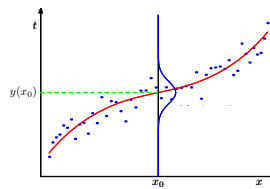
12

Explicitly Modeling Uncertainty: The Generative Picture



The generating process...

$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$

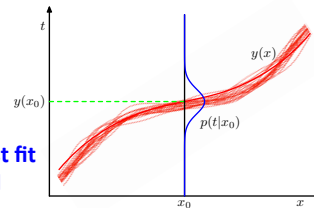


... generates data ...

Each curve here is a best fit to one of the generated datasets.

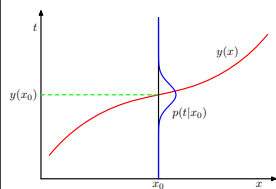
Each is an estimate of the true generating process, as mediated by the dataset

Suppose we generate 10,000 datasets (each slightly different) but from the same generating process...



13

Explicitly Modeling Uncertainty: The Generative Picture

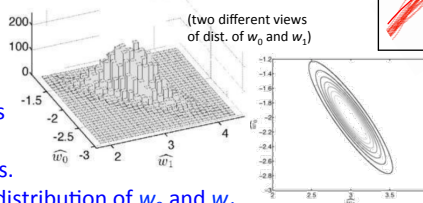


The generating process...

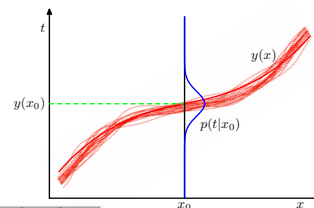
$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$

We can estimate the uncertainty in parameter values by looking at the empirical distribution of different parameter values we get from each of the different dataset estimates.

E.g., here is the empirical distribution of w_0 and w_1 ...

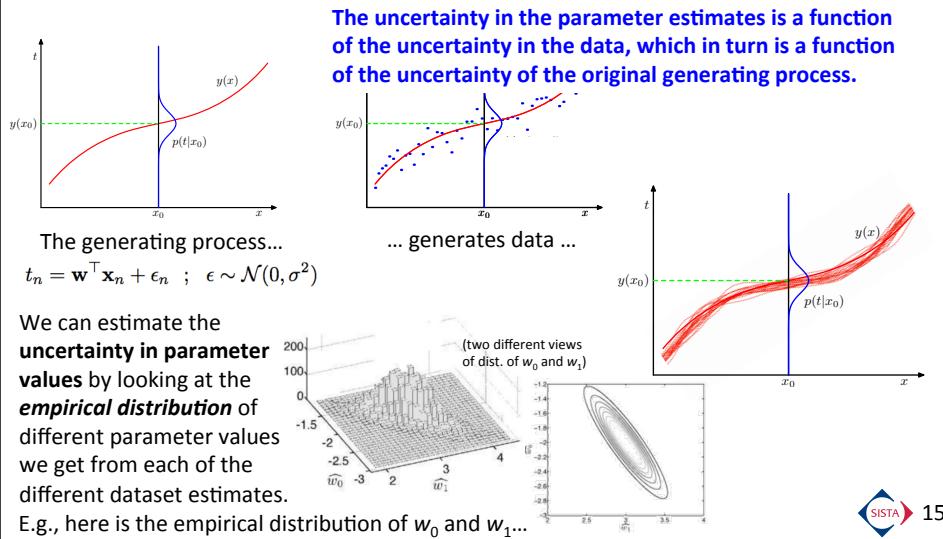


Suppose we generate 10,000 datasets (each slightly different) but from the same generating process...

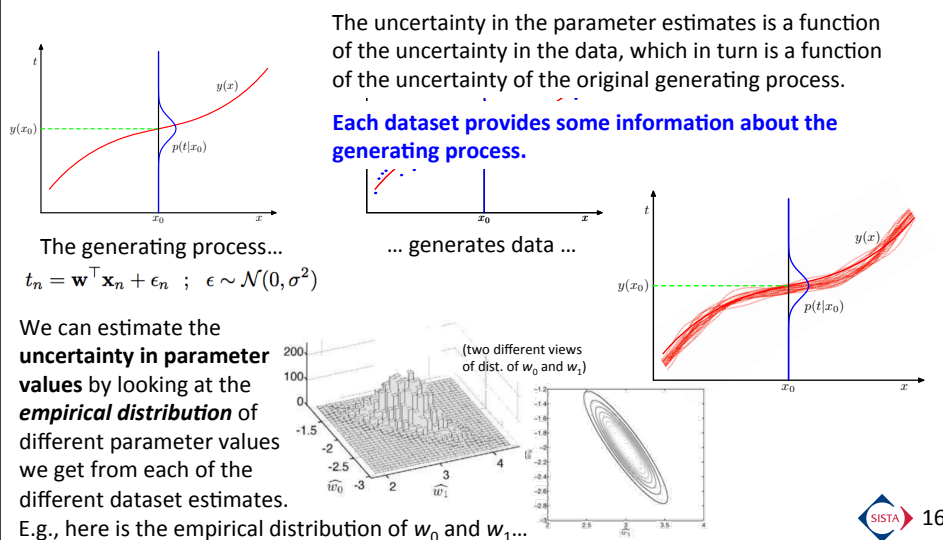


14

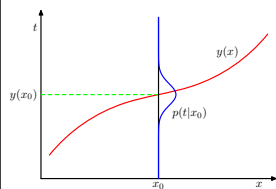
Explicitly Modeling Uncertainty: The Generative Picture



Explicitly Modeling Uncertainty: The Generative Picture



Explicitly Modeling Uncertainty: The Generative Picture

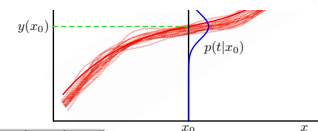
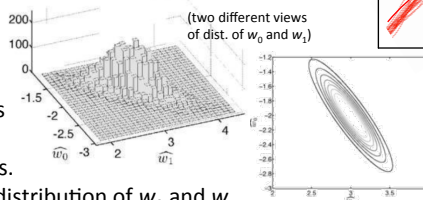


The generating process...

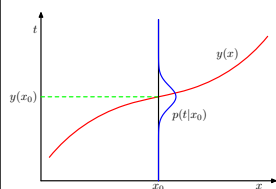
$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$

We can estimate the **uncertainty in parameter values** by looking at the **empirical distribution** of different parameter values we get from each of the different dataset estimates.

E.g., here is the empirical distribution of w_0 and w_1 ...



Explicitly Modeling Uncertainty: The Generative Picture

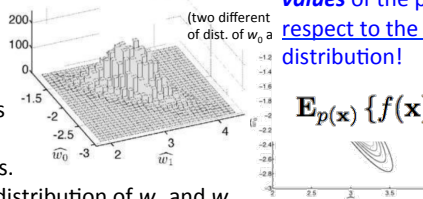


The generating process...

$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$

We can estimate the **uncertainty in parameter values** by looking at the **empirical distribution** of different parameter values we get from each of the different dataset estimates.

E.g., here is the empirical distribution of w_0 and w_1 ...



Can we estimate the **uncertainty in our parameters** from a single dataset?

YES! ... by considering the **expected values** of the parameters with **respect to the generating process distribution!**

$$\mathbf{E}_{p(\mathbf{x})} \{f(\mathbf{x})\} = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$



Side Note

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n|\mathbf{x}_n \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$$

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$$

In the following, easier to deal with multivariate Gaussian rather than product of individual Gaussians



Deriving the Expectation of $\hat{\mathbf{w}}$ w.r.t. the **generating distribution**

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = \int \hat{\mathbf{w}} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) d\mathbf{t}$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} \quad \text{Substitute in } \hat{\mathbf{w}}$$

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \int \mathbf{t} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) d\mathbf{t}$$

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t}\} \quad p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$$

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{w} \quad \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t}\} = \mathbf{X} \mathbf{w}$$

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = \mathbf{w}$$

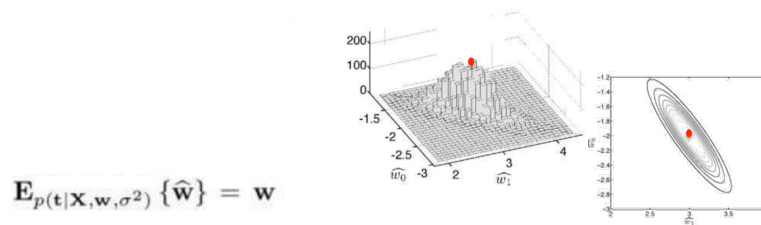
Our estimate $\hat{\mathbf{w}}$ of \mathbf{w} is inherently **unbiased**!

This also means any variance in the estimate is encapsulated in its **(co)variance (matrix) of $\hat{\mathbf{w}}$**

Deriving the Expectation of $\hat{\mathbf{w}}$ w.r.t. the generating distribution

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = \int \hat{\mathbf{w}} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) d\mathbf{t}$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} \quad \text{Substitute in } \hat{\mathbf{w}}$$



Our estimate $\hat{\mathbf{w}}$ of \mathbf{w} is inherently **unbiased**!

This also means any variance in the estimate is encapsulated in its **(co)variance (matrix) of $\hat{\mathbf{w}}$**

Let's be clear on what we mean by "Maximum likelihood of the mean of a linear Gaussian model is unbiased"

- Here we mean: if we take repeated samples of size N from a Gaussian generator distribution with true parameters \mathbf{w} , then the collection of maximum likelihood estimated $\hat{\mathbf{w}}$'s will be "centered" (in the sense of the arithmetic mean) around the true \mathbf{w} .
- Not to be confused with:
 - "better" estimation by collecting more data (it is true that more data makes better estimation here, but that's a sample size effect)
 - "UMVU" estimator: *Uniformly minimum variance unbiased estimator*. This is a stronger claim: an UMVU estimator is also an unbiased estimator that is "closest" to the true parameters, for a given sample size.
- And, there could be:
 - Other *biased* estimators that tend to be closer (from less data)

Derive the Covariance of $\hat{\mathbf{w}}$ w.r.t. the generating distribution

$$\text{cov}\{\mathbf{x}\} = \mathbb{E}_{P(\mathbf{x})} \left\{ \mathbf{x} \mathbf{x}^\top \right\} - \mathbb{E}_{P(\mathbf{x})} \{\mathbf{x}\} \mathbb{E}_{P(\mathbf{x})} \{\mathbf{x}\}^\top$$

$$\begin{aligned} \text{cov}\{\hat{\mathbf{w}}\} &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \right\} - \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\}^\top \\ &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \right\} - \mathbf{w} \mathbf{w}^\top \end{aligned} \quad \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = \mathbf{w}$$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \right\} &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}) ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t})^\top \right\} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{t} \mathbf{t}^\top \right\} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}. \end{aligned}$$

$$\text{cov}\{\mathbf{t}\} = \sigma^2 \mathbf{I} = \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{t} \mathbf{t}^\top \right\} - \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t}\} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t}\}^\top$$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{t} \mathbf{t}^\top \right\} &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t}\} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t}\}^\top + \sigma^2 \mathbf{I} \\ &= \mathbf{X} \mathbf{w} (\mathbf{X} \mathbf{w})^\top + \sigma^2 \mathbf{I} \\ &= \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top + \sigma^2 \mathbf{I}. \end{aligned}$$

$$\begin{aligned} \text{cov}\{\hat{\mathbf{w}}\} &= \mathbf{w} \mathbf{w}^\top + \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} - \mathbf{w} \mathbf{w}^\top \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \right\} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &\quad + \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \mathbf{w} \mathbf{w}^\top + \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}. \end{aligned}$$

$$\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} = -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} \quad \leftarrow \text{Remember!}$$

$$\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = - \left(\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right)^{-1} \quad \blacktriangleright 23$$

The Fisher Information

$$\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = - \left(\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right)^{-1}$$

$$\mathcal{I} = \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ - \frac{\partial^2 \log p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right\}$$

$$\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} = -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

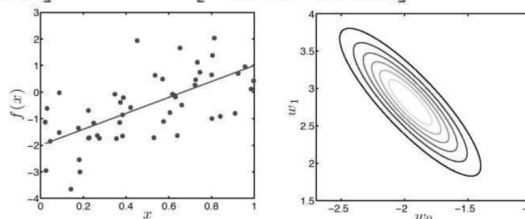
$$\mathcal{I} = \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} \right\}$$

$$\mathcal{I} = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}.$$

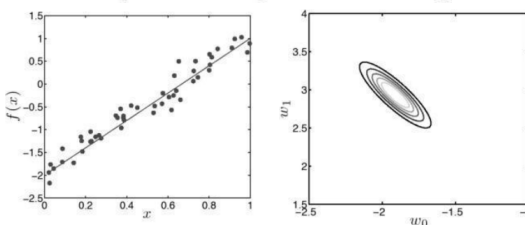
The elements of \mathcal{I} tell us how much information the data provides about the particular parameter (diagonal elements) or pairs of parameters (off-diagonal elements).

Example

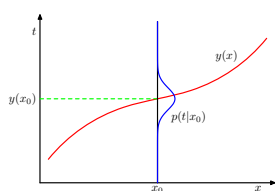
$$\mathcal{I} = \begin{bmatrix} 50.0000 & 24.3311 \\ 24.3311 & 15.8953 \end{bmatrix}, \text{cov}\{\hat{\mathbf{w}}\} = \begin{bmatrix} 0.0784 & -0.1200 \\ -0.1200 & 0.2466 \end{bmatrix}$$



$$\mathcal{I} = \begin{bmatrix} 1.2500 \times 10^3 & 0.6083 \times 10^3 \\ 0.6083 \times 10^3 & 0.3974 \times 10^3 \end{bmatrix}, \text{cov}\{\hat{\mathbf{w}}\} = \begin{bmatrix} 0.0031 & -0.0048 \\ -0.0048 & 0.0099 \end{bmatrix}$$

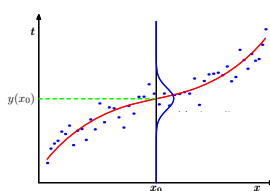


Back to: The Generative Picture

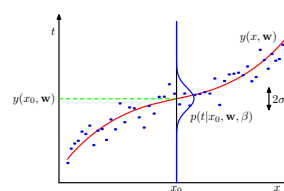


The generating process...

$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n; \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$



... generates data ...



... that we fit a model to

How about modeling the uncertainty in our predictions?

Want model of region of values in which our prediction might fall

estimated parameters

$$p(\hat{\mathbf{t}}|\mathbf{X}, \hat{\mathbf{w}}, \hat{\sigma}^2) = \prod_{n=1}^N p(\hat{t}_n|\mathbf{x}_n, \hat{\mathbf{w}}, \hat{\sigma}^2)$$

predictive distribution

$$= \prod_{n=1}^N \mathcal{N}(\hat{\mathbf{w}}^\top \mathbf{x}_n, \hat{\sigma}^2)$$