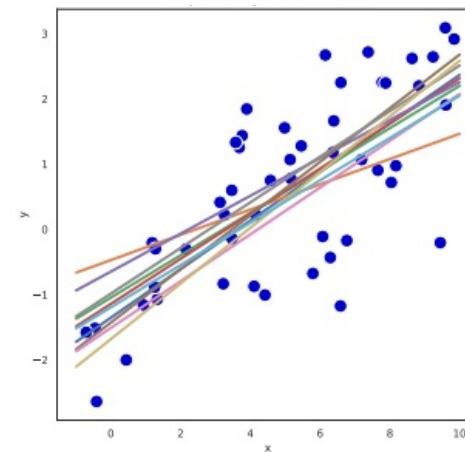
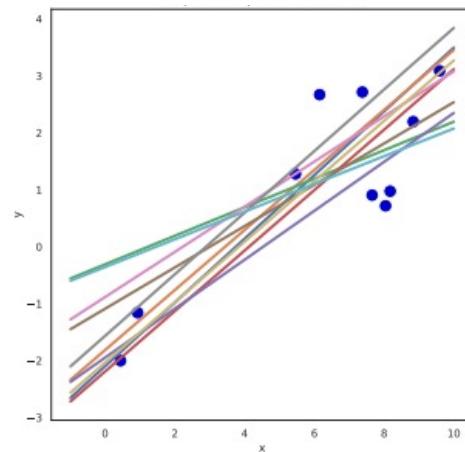
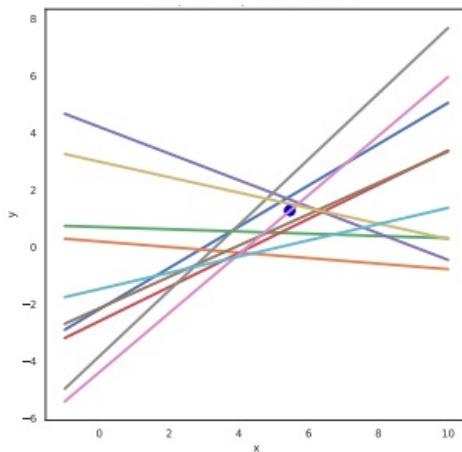


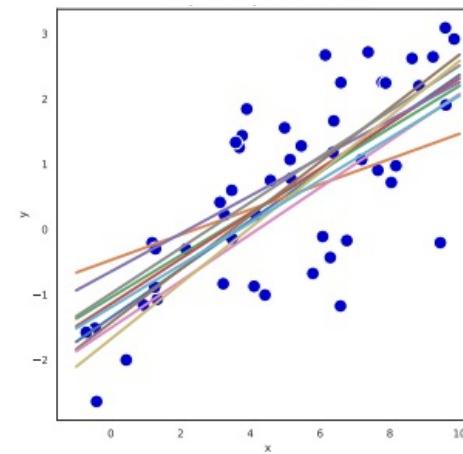
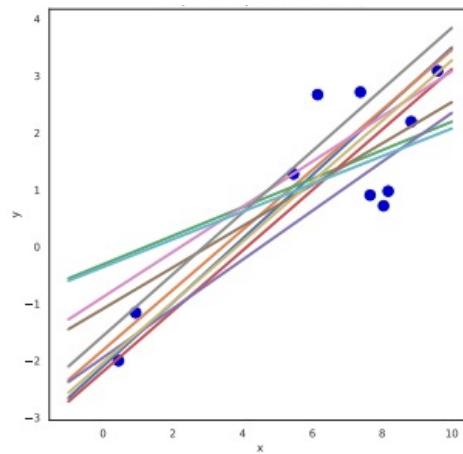
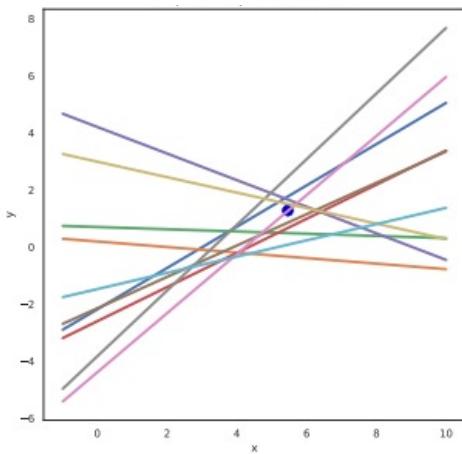
Machine Learning

Bayesian Regression

Jian Liu



Noisy Targets – Target Distribution



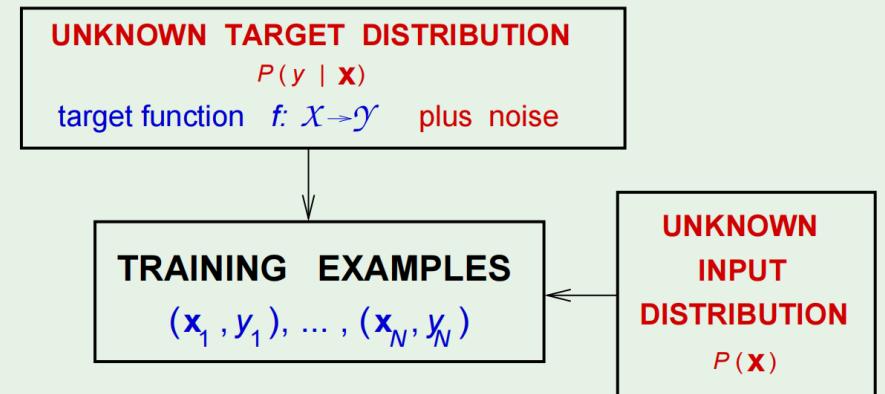
Noisy Targets – Target Distribution

Both convey probabilistic aspects of \mathbf{x} and y

The target distribution $P(y | \mathbf{x})$
is what we are trying to learn

The input distribution $P(\mathbf{x})$
quantifies relative importance of \mathbf{x}

Merging $P(\mathbf{x})P(y|\mathbf{x})$ as $P(\mathbf{x}, y)$
mixes the two concepts



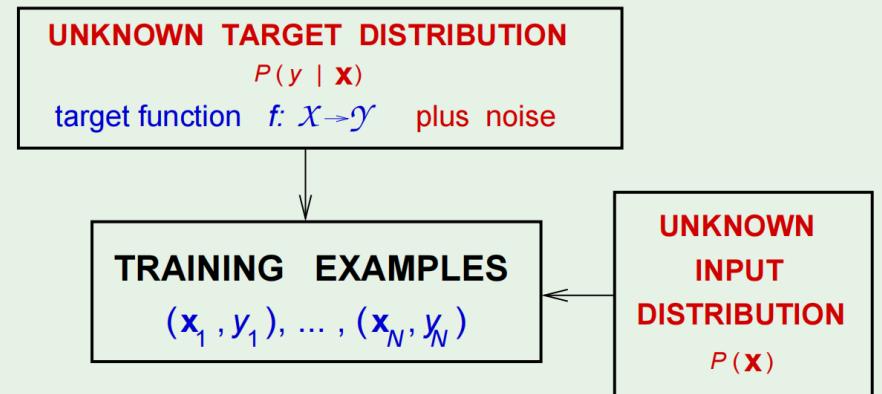
Probabilistic Machine Learning

Both convey probabilistic aspects of \mathbf{x} and y

The target distribution $P(y | \mathbf{x})$
is what we are trying to learn

The input distribution $P(\mathbf{x})$
quantifies relative importance of \mathbf{x}

Merging $P(\mathbf{x})P(y|\mathbf{x})$ as $P(\mathbf{x}, y)$
mixes the two concepts



Probabilistic Machine Learning

- Key difference in Frequentist vs Bayesian paradigms:
 - **Frequentist:** model parameters θ are **fixed**; computed using some *estimator* such as **MLE**
 - Confidence in estimates for θ evaluated through multiple experiments (**cross-validation**) for different data sets
 - **Bayesian:** model parameters θ are **random variables**
 - Here, there is only **one** data set (actually observed) and **uncertainty in θ** is expressed as a probability distribution over θ

Probabilistic Machine Learning

Bayes' Law

- Using the sum and product rules of probability for discrete values of θ (or events):

$$P(\theta | X) = \frac{P(X | \theta)P(\theta)}{P(X)}$$

where, the marginal probability $P(X) = \sum_{\theta} P(X | \theta)P(\theta)$

- For continuous variables, i.e. all possible values of θ we can do the same:

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

except now, $p(\theta)$ is a **continuous prior distribution** rather than **discrete** probabilities

- Now the marginal probability density function $p(X) = \int p(X|\theta)p(\theta)d\theta$
- So we can say: *posterior \propto likelihood \times prior*

Linear Regression: OLS

$$\hat{y}_i(x_i, \omega) = \sum_{j=0}^M \omega_j \phi_j(x_i) = \omega^T \phi(x_i)$$

- By defining design matrix $\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_M(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \cdots & \phi_M(x_2) \\ \vdots & \vdots & & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \cdots & \phi_M(x_N) \end{pmatrix}$
- ω_{OLS} is given by solving: $\Phi^T \mathbf{y} = \Phi^T \Phi \omega$ **The normal equation**

$$\omega_{OLS} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

Linear Regression: Probabilistic View

- Consider again the linear relationship between target y & predictor x :

$$\mathbf{y} = \hat{y}(\mathbf{x}, \boldsymbol{\omega}) + \epsilon$$

where, ϵ is some *unexplained* noise; $\hat{y}(\cdot)$ is a linear combination of basis functions

- Lets assume ϵ is a univariate Gaussian variable with zero mean and precision $\beta = \frac{1}{\sigma^2}$
 - β used for convenience of notation
- Then the probability distribution of a target y *conditioned* on $\mathbf{x}, \boldsymbol{\omega}$ is given by:

$$p(y | \mathbf{x}, \boldsymbol{\omega}, \beta) = \mathcal{N}(y | \hat{y}(\mathbf{x}, \boldsymbol{\omega}), \beta^{-1})$$

MLE (Maximum Likelihood Estimation)

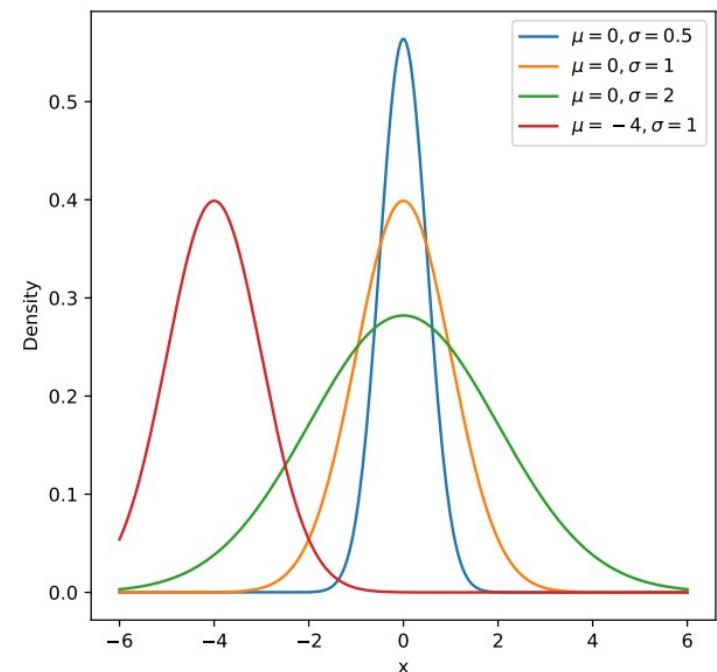
Recap: Univariate Gaussian distribution

- Occurs often in Nature: e.g. height, IQ
- Fixed length sums of variables of any distribution follow approximate Gaussian: **Central Limit Theorem**
- μ : position of centre; σ : width of distribution

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \quad (1)$$

$$\mathbb{E}[\mathcal{N}(x | \mu, \sigma^2)] = \mu \quad (2)$$

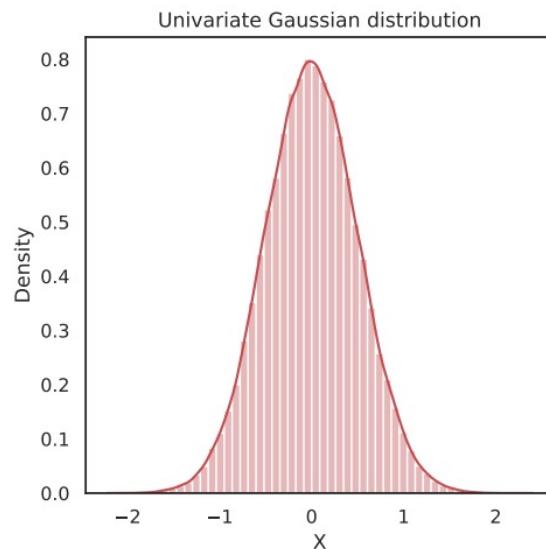
$$\text{var}[\mathcal{N}(x | \mu, \sigma^2)] = \sigma^2$$



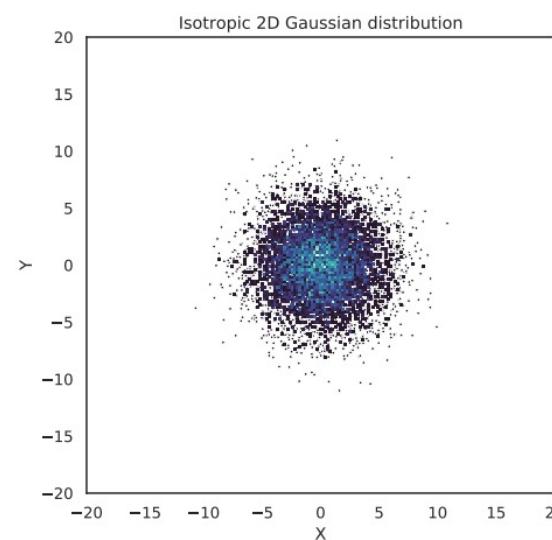
MLE (Maximum Likelihood Estimation)

Recap: Multivariate Gaussian Distribution

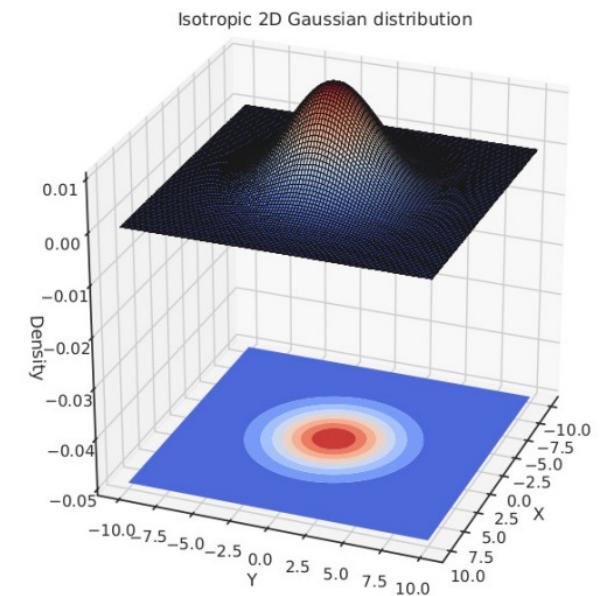
- Ubiquity of the Gaussian distribution can be explained by the **Central Limit Theorem**
- This holds for Univariate and Multivariate Gaussian distributions
- What is a **multivariate** distribution?



(a) Univariate Gaussian distribution



(b) Samples from bivariate Gaussian



(c) Density of bivariate Gaussian

MLE (Maximum Likelihood Estimation)

Recap: Multivariate Gaussian Distribution

- For an N -dim vector \mathbf{x} the multivariate Gaussian distribution is given by:

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{N/2}} \frac{1}{|\det(\boldsymbol{\Sigma})|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

- Here, $\boldsymbol{\mu}$ is an N -dim mean vector and $\boldsymbol{\Sigma}$ is an $N \times N$ covariance matrix
- For example for a bivariate (2D) Gaussian: $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{xx}^2 & \sigma_{xy} \\ \sigma_{yx} & \sigma_{yy}^2 \end{bmatrix}$

MLE (Maximum Likelihood Estimation)

- Parameters to estimate when fitting a multivariate Gaussian to data are: μ, Σ
- Lets look at deriving MLE for μ, Σ
- Given a dataset $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ of N i.i.d data points each forming a row of the matrix,
- Using the product rule we define the likelihood function as:

$$p(\mathbf{X}|\mu, \Sigma) = \prod_{i=1}^N p(\mathbf{x}_i|\mu, \Sigma)$$

- And we know that for a single data point \mathbf{x}_i we have:

$$p(\mathbf{x}_i|\mu, \Sigma) = \mathcal{N}(\mathbf{x}_i | \mu, \Sigma) = \frac{1}{(2\pi)^{N/2}} \frac{1}{|\det(\Sigma)|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu) \right\}$$

MLE (Maximum Likelihood Estimation)

- To derive MLE for μ & Σ we need to maximise the likelihood function w.r.t each parameter
- I.e. we need to solve:

$$\frac{\partial}{\partial \mu} p(\mathbf{X} | \mu, \Sigma) |_{\mu=\mu_{ML}} = 0$$

$$\frac{\partial}{\partial \Sigma} p(\mathbf{X} | \mu, \Sigma) |_{\Sigma=\Sigma_{ML}} = 0$$

- Easier to work with the log-likelihood function (\mathcal{L}):

$$\mathcal{L} = \ln p(\mathbf{X} | \mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln \det(\Sigma) - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu)$$

MLE (Maximum Likelihood Estimation)

- Considering just the terms in (\mathcal{L}) dependent on μ & Σ :

$$\mathcal{L}(\mu) = \sum_{i=1}^N (\mathbf{x}_i - \mu)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \mu)$$

$$\mathcal{L}(\Sigma) = -\frac{N}{2} \ln \det(\Sigma) - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mu)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \mu)$$

- Solving $\frac{\partial \mathcal{L}(\mu)}{\partial \mu} = 0$ & $\frac{\partial \mathcal{L}(\Sigma)}{\partial \Sigma} = 0$ gives:

$$\mu_{ML} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

$$\boldsymbol{\Sigma}_{ML} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^T$$

Linear Regression: MLE & OLS

- Given a data set $\mathbf{X} = \{\mathbf{x}_i\}_{i=1\dots N}$ & targets $\mathbf{y} = \{y_i\}_{i=1\dots N}$, assuming data are *i.i.d*,
- We can formulate the likelihood function as a function of ω, β :

$$p(\mathbf{y} | \mathbf{X}, \omega, \beta) = \prod_{i=1}^N \mathcal{N}(y_i | \omega^T \phi(x_i), \beta^{-1})$$

where we have used the previous result: $\hat{y}(x_i, \omega) = \omega^T \phi(x_i)$

- But we have seen this form before for **MLE of multivariate Gaussians...**

we can express the log-likelihood \mathcal{L} as:

$$\mathcal{L} = \ln p(\mathbf{y} | \mathbf{X}, \omega, \beta) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln \beta - \frac{1}{2} \sum_{i=1}^N \beta(y_i - \omega^T \phi(x_i))^2$$

Linear Regression: MLE & OLS

- Looking at the quadratic term in \mathcal{L} :

$$\mathcal{L} = \ln p(\mathbf{y} | \mathbf{X}, \omega, \beta) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln \beta - \frac{\beta}{2} R(\omega)$$

$$R(\omega) = \sum_{i=1}^N (y_i - \omega^T \phi(x_i))^2 = \sum_{i=1}^N (y_i - \hat{y}_i(x_i, \omega))^2$$

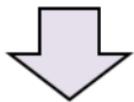
- But this is the **same** as the **sum of squared residuals** in OLS!
- ω which **maximises** \mathcal{L} also **minimises** R : $\arg \max_{\omega} \mathcal{L} \equiv \arg \min_{\omega} (-\mathcal{L})$
$$\arg \min_{\omega} (-\mathcal{L}(\omega)) = \arg \min_{\omega} \sum_{i=1}^N (y_i - \hat{y}_i(x_i, \omega))^2 = \arg \min_{\omega} R(\omega)$$
- So $\omega_{MLE} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$ under **additive Gaussian noise assumption**
 $\omega_{OLS} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$ **MLE = OLS**

Linear Regression: MLE & OLS

Regression

Optimization:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \quad \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2$$



Solution:

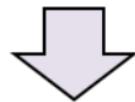
$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad \equiv$$

Assumption: None

Maximum-Likelihood

Optimization:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \quad \left(\frac{1}{\sqrt{(2\pi\sigma^2)^d}} \right)^N \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 \right\}$$



Solution:

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

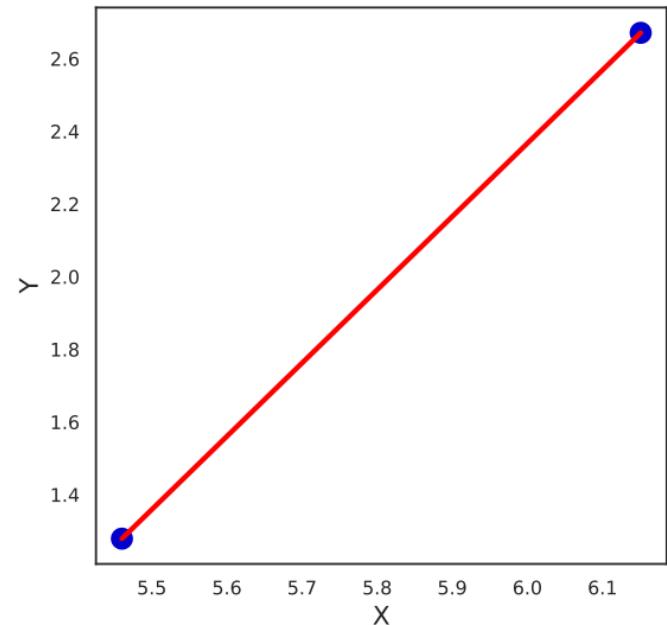
Assumption: $\mathbf{y} - \mathbf{X}\boldsymbol{\theta} \sim \text{Gaussian}(0, \sigma^2 \mathbf{I})$

Linear Regression: Bayesian perspective

- Consider linear regression with a polynomial of degree = 1; $\hat{y}_i = \omega_0 + \omega_1 x_i$
- If we have just two predictors and their targets: $\{x_i, y_i\}_{i=1,2}$, we can fit a line exactly

- Here we have two unknowns and two observations
- Can determine the slope (ω_1) and intercept (ω_0) just using the observations
- We can solve for ω_1 as follows and then estimate ω_0 :

$$\omega_1 = \frac{y_2 - y_1}{x_2 - x_1} \quad (13)$$



Linear Regression: Bayesian perspective

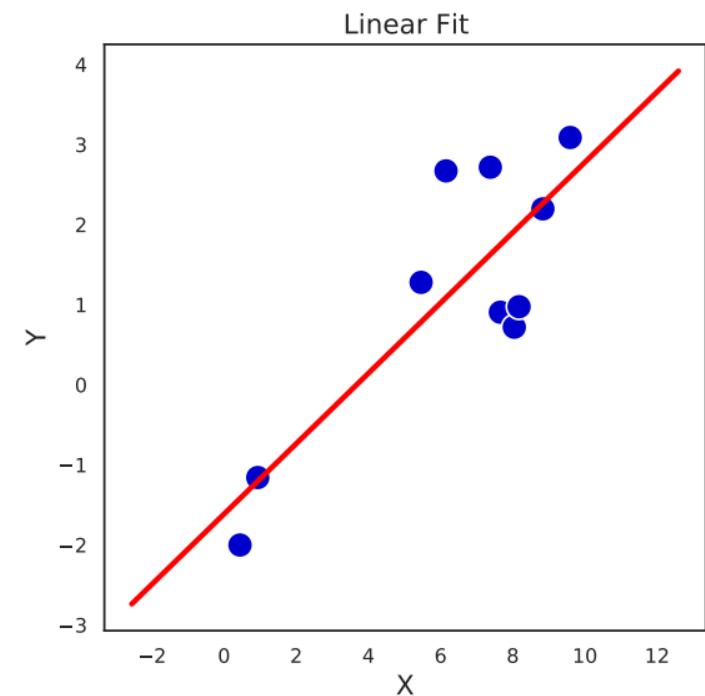
- But what if we had more than two points?

- Then we have an *overdetermined* system
- We can use OLS or assume a noise model
- Then system of linear equations are:

$$y_i = \omega_0 + \omega_1 x_i + \epsilon_i \quad (14)$$

where, we can assume $\epsilon \sim \mathcal{N}(0, \sigma^2)$

- And this leads to $\omega_{MLE} \equiv \omega_{OLS}$



Linear Regression: Bayesian perspective

- What if we have an *underdetermined* system? E.g. **just one observation**
- We have **two unknowns** and **one observation**: $y_1 = \omega_0 + \omega_1 x_1 + \epsilon_1$
- What can we do if we want to fit a line to this?

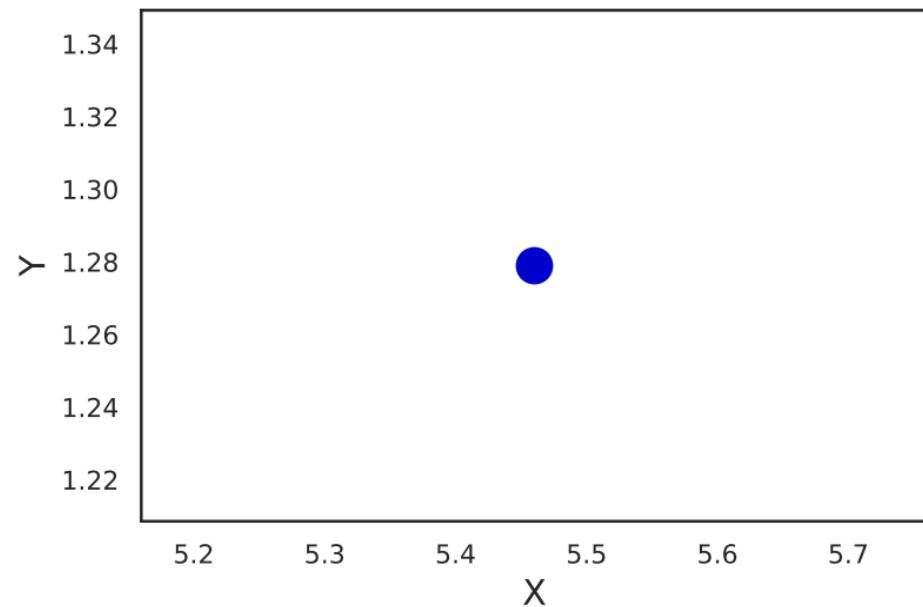
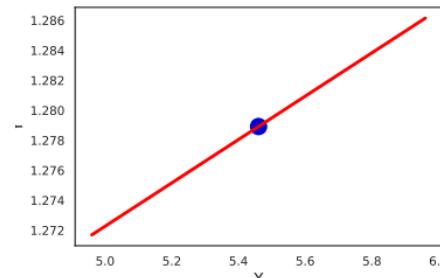


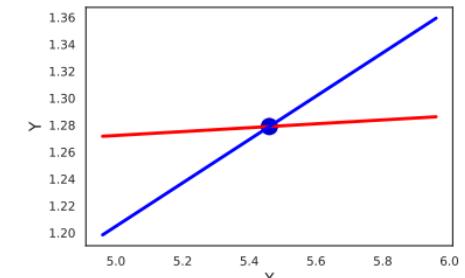
Figure: Just one observations (x_1, y_1)

Linear Regression: Bayesian perspective

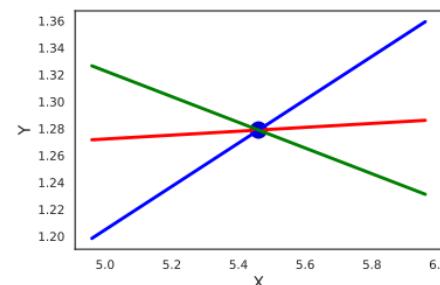
- What if we have an *underdetermined* system? E.g. just one observation
- We have two unknowns and one observation: $y_1 = \omega_0 + \omega_1 x_1 + \epsilon_1$
- We have an infinite number of possible solutions..
- How can we represent a **family** of ω_0 ?



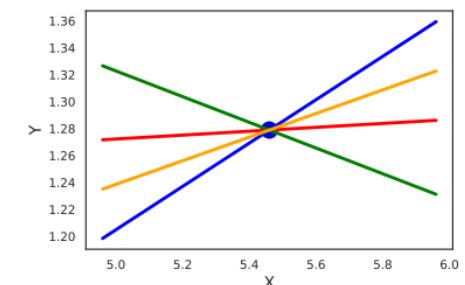
(a)



(b)



(c)



(d)

Figure: (a): $\omega_0 = 1.2$, (b): $\omega_0 = 0.4$, (c): $\omega_0 = 1.8$,
(d): $\omega_0 = 0.8$

Linear Regression: Bayesian perspective

- We can make an **assumption for the distribution** of the unknown parameter ω_0 !
- And that is **Bayesian Inference**!
- Given a data set $\mathbf{X} = \{\mathbf{x}_i\}_{i=1\dots N}$ & targets $\mathbf{y} = \{y_i\}_{i=1\dots N}$, assuming data are *i.i.d.*,
- As previously the likelihood function is given by:

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\omega}, \beta) = \prod_{i=1}^N \mathcal{N}(y_i | \boldsymbol{\omega}^T \phi(\mathbf{x}_i), \sigma^2)$$

- We tried to **maximise the likelihood** with respect to $\boldsymbol{\omega}$
- Bayesian perspective: based on some **a priori belief** about $\boldsymbol{\omega}$, integrate across $\boldsymbol{\omega}$
- I.e. instead of fixed $\boldsymbol{\omega}$ we look at expectation of likelihood for a range of plausible $\boldsymbol{\omega}$

Linear Regression: Bayesian perspective

- For Bayesian inference we first define a **prior distribution** on the unknowns/parameters (ω)
- Prior represents our belief about ω **before we see the data**
- Linear regression: let's consider a Gaussian prior on the intercept (ω_0) (ignore ω_1 for now):

$$p(\omega_0) \sim \mathcal{N}(0, \alpha)$$

- Given N i.i.d observations $\{x_i, y_i\}_{i=1\dots N}$ and $p(\omega_0)$ the posterior distribution of ω_0 is:

$$p(\omega_0 | \mathbf{y}, \mathbf{x}, \omega_1, \sigma^2) \propto \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \omega_1 x_i - \omega_0)^2 - \frac{\omega_0^2}{2\alpha} \right\}$$

$$\ln p(\omega_0 | \mathbf{y}, \mathbf{x}, \omega_1, \sigma^2) = \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \omega_1 x_i - \omega_0)^2 - \frac{\omega_0^2}{2\alpha} \right\} + const$$

Linear Regression: Bayesian perspective

- Once again by **completing the square** we find:

$$\ln p(\omega_0 \mid \mathbf{y}, \mathbf{x}, \omega_1, \sigma^2) = -\frac{1}{\sigma_{post}^2} (\omega_0 - \mu_{post})^2 + const$$

$$p(\omega_0 \mid \mathbf{y}, \mathbf{x}, \omega_1, \sigma^2) \sim \mathcal{N}(\mu_{post}, \sigma_{post}^2)$$

$$\mu_{post} = \frac{N\alpha}{N\alpha + \sigma^2} \mu_{ML}, \quad \sigma_{post}^2 = \frac{\sigma^2 \alpha}{N\alpha + \sigma^2}$$

- Where μ_{ML} is the **MLE for the intercept ω_0** (refer to C. Bishop's book pg. 142)
- For MLE we estimated all regression weights jointly
- We can infer the **posterior distribution over all weights** by assuming a **prior $p(\omega)$** over all

Linear Regression: Bayesian perspective

- Previously, we assumed a prior over the intercept for a linear regression function
- What about a **prior over all** regression weights?
- Given a data set N i.i.d observations: $\mathbf{y} = \hat{y}(\mathbf{x}, \boldsymbol{\omega}) + \boldsymbol{\epsilon}$

$$p(\mathbf{y} | \mathbf{x}, \boldsymbol{\omega}, \beta) = \prod_{i=1}^N \mathcal{N}(y_i | \boldsymbol{\omega}^T \phi(x_i), \beta^{-1})$$
$$p(\boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\omega} | \mathbf{0}, \alpha^{-1} \mathbb{I})$$

- Then posterior distribution over weights is given by:

$$p(\boldsymbol{\omega} | \mathbf{y}, \mathbf{x}, \beta) \propto p(\mathbf{y} | \mathbf{x}, \boldsymbol{\omega}, \beta) p(\boldsymbol{\omega})$$

- We will treat the precision β as a **known constant** - but this may be inferred as well

Linear Regression: Bayesian perspective

- We know if $p(\omega)$ is a **conjugate prior** then $p(\omega | \mathbf{y}, \mathbf{x})$ is a Gaussian distribution
- As seen previously, given $p(\omega) = \mathcal{N}(\omega | \mathbf{0}, \alpha^{-1}\mathbb{I})$:

$$\ln p(\mathbf{w} | \mathbf{y}, \mathbf{x}) = -\frac{\beta}{2} \sum_{i=1}^N \left\{ y_i - \mathbf{w}^T \phi(x_i) \right\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const}$$

- Once again by *completing the square*:

$$p(\omega | \mathbf{y}, \mathbf{x}) = \mathcal{N}(\omega | \mathbf{m}_{post}, \mathbf{S}_{post})$$

$$\mathbf{m}_{post} = \beta \mathbf{S}_{post} \Phi^T \mathbf{y}, \quad \mathbf{S}_{post}^{-1} = \alpha \mathbb{I} + \beta \Phi^T \Phi$$

where, Φ is the design matrix, $\mathbf{m}_{post}, \mathbf{S}_{post}$ are the mean and covariance of the posterior distribution for ω

Linear Regression: Bayesian perspective

- MLE for linear regression: $\omega_{MLE} = \Phi^\dagger$
- MLE gives point estimates for ω , but a predictive distribution in *data space*
- Remember:
 $p(y_i | x_i, \omega, \beta) = \mathcal{N}(y_i | \omega^T \phi(x_i), \beta^{-1})$
- Ex: Fitting degree=1 polynomial; $\beta = 1.0$
- Predictive distribution is:
 $\mathcal{N}(y | (\omega_0 + \omega_1 x), 1.0)$
- What do you notice about the MLE distribution?

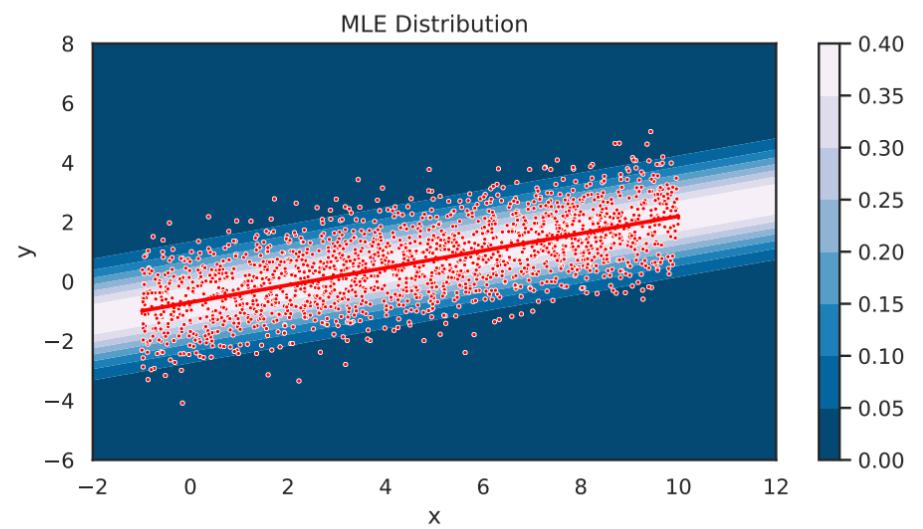


Figure: Polynomial degree=1 fit to data; contour plot depicts predictive distribution

Linear Regression: Bayesian perspective

- Bayesian inference for linear regression: posterior distribution of weights $p(\omega | \mathbf{y}, \mathbf{x})$
- I.e. we have a predictive distribution for each **possible ω** from $p(\omega | \mathbf{y}, \mathbf{x})$
- Consider $N = 10$ i.i.d observations, $\beta = 1.0$, prior $p(\omega) \sim \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbb{I})$
- We infer $p(\omega | \mathbf{y}, \mathbf{x}) = \mathcal{N}(\omega | \mathbf{m}_{post}, \mathbf{S}_{post})$: What do you **notice**?

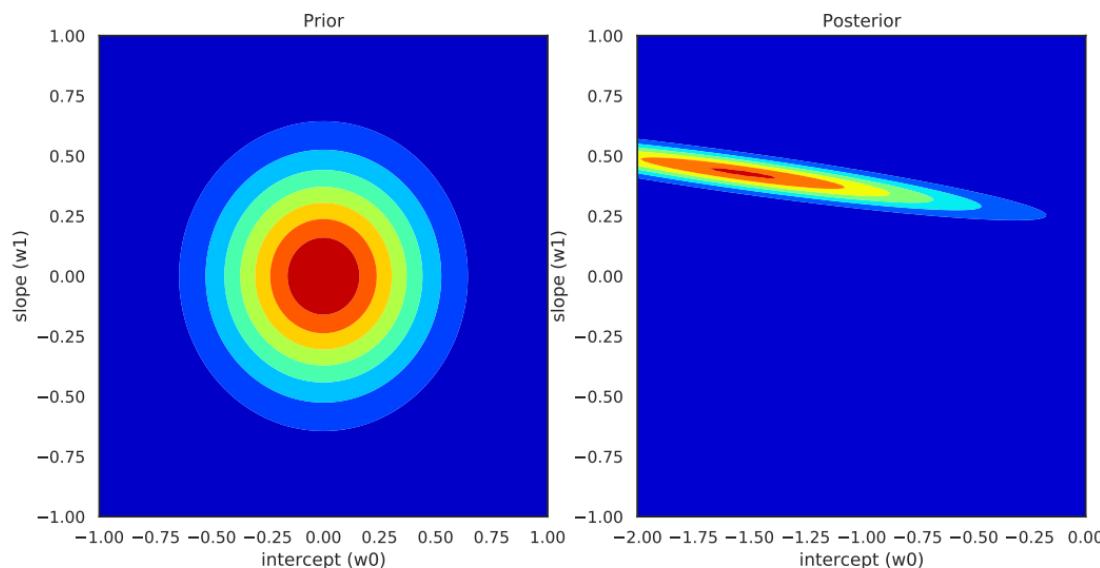
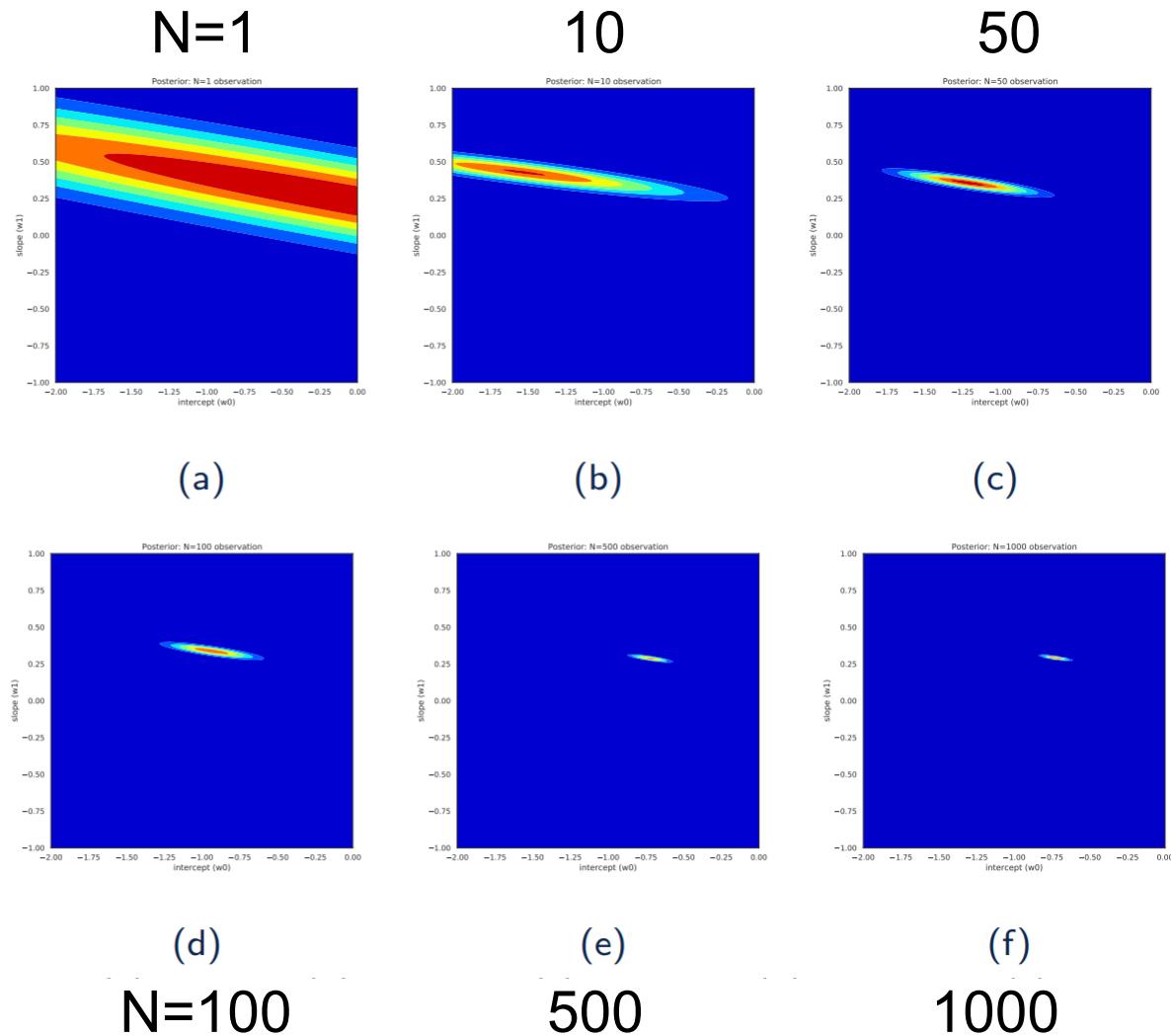


Figure: Left: Prior distribution, $\alpha = 10.0$; Right: Posterior distribution after $N = 10$ observations

Linear Regression: Bayesian perspective

- Let's see what the impact of varying the number of observations is on the posterior

- What do you notice?
- x-axis: intercept (ω_0);
y-axis: slope (ω_1)



Linear Regression: Bayesian perspective

- *Model uncertainty reduces with increase in observations*
- What about the predictive distribution?
- Let's sample 10 instances of ω from $\mathcal{N}(\mathbf{m}_{post}, \mathbf{S}_{post})$
- What do you notice?

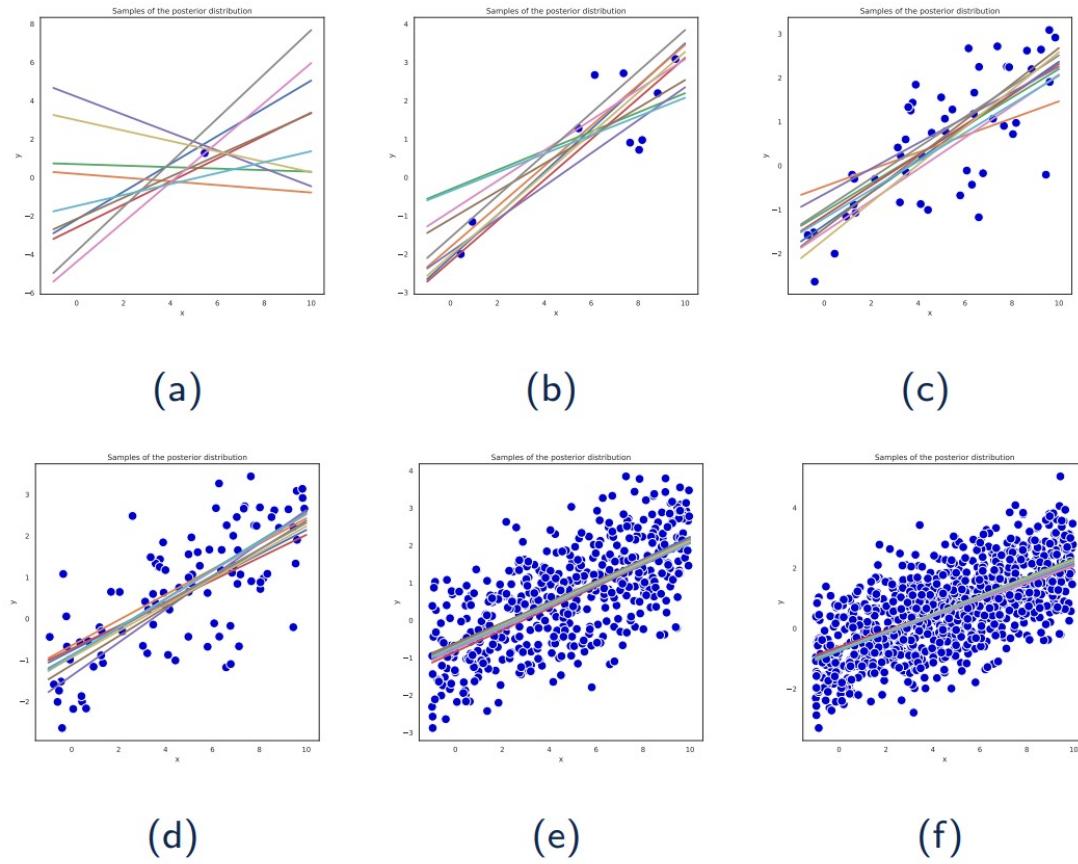


Figure: (a) $N = 1$; (b) $N = 10$; (c) $N = 50$; (d) $N = 100$; (e) $N = 500$; (f) $N = 1000$

Summary

- OLS = MLE (Point estimation)
- Bayesian linear regression models
(Distribution of model parameters)
- Data are too few or too noisy