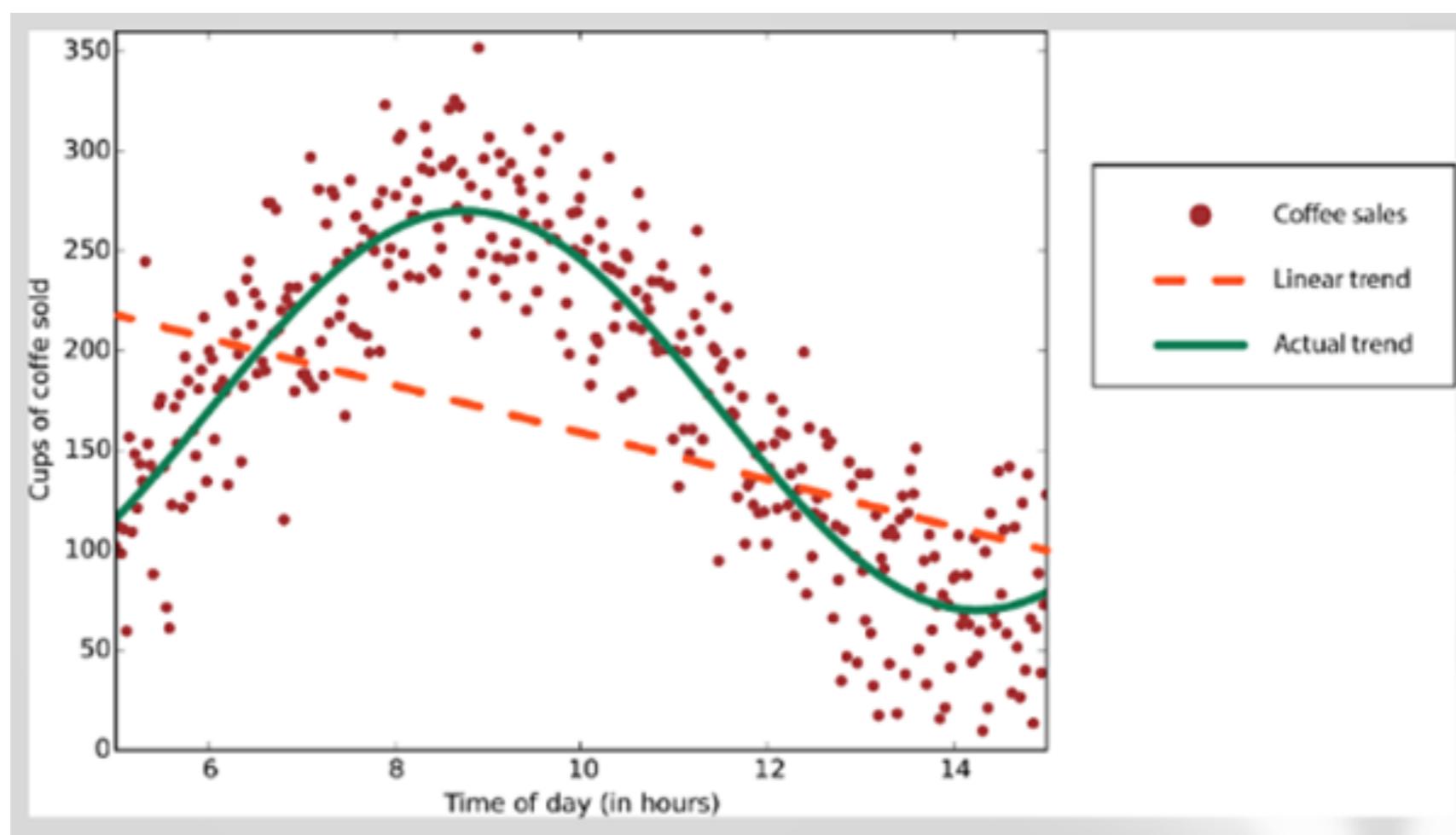


LINEAR REGRESSION

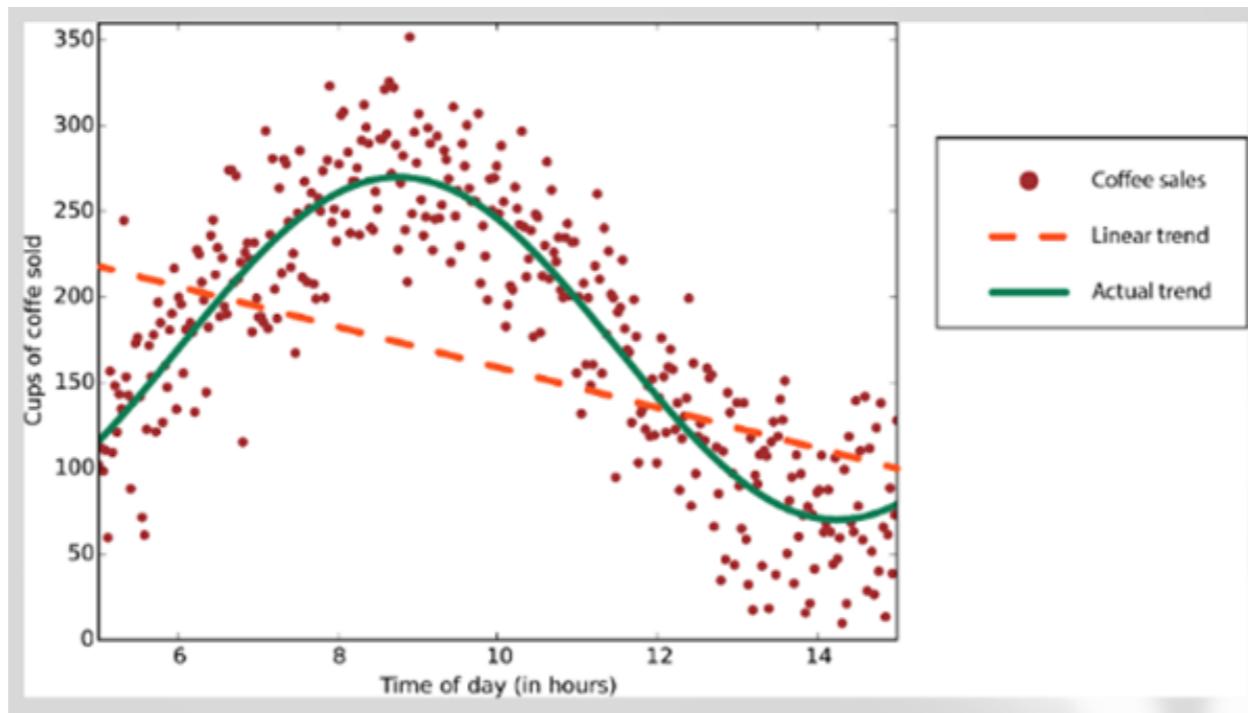
Dr. Srijith P K



Supervised learning : Regression



Supervised learning : Regression



Real valued targets (outputs)

Generalization performance

Goal is to learn a function which maps inputs to outputs so that it will predict well on future data points

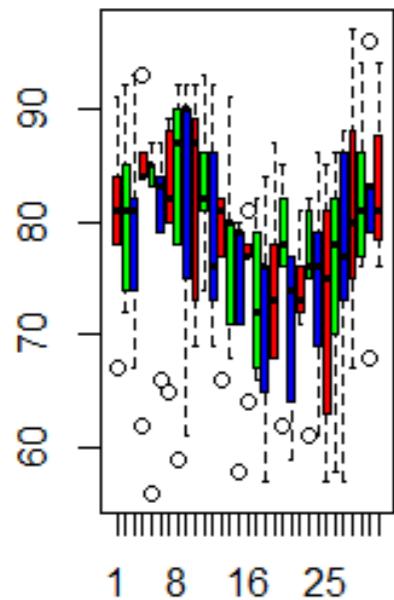
Airquality data.

- Data set has various air quality parameters in New York city.
- These are the parameters in the data set:
 - Daily temperature from May to August
 - Solar radiation data
 - Ozone data
 - Wind data
- Goal : predict the temperature for a particular month in New York using solar radiation, ozone and wind data.

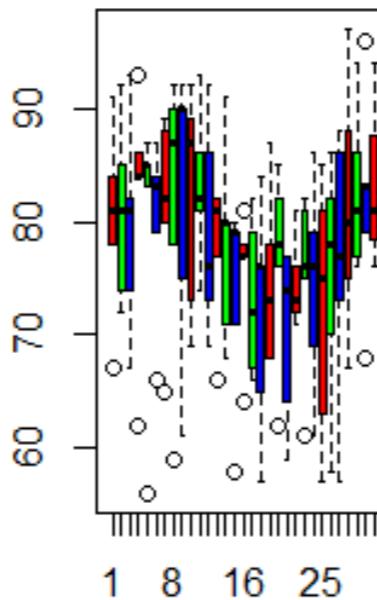
Airquality data

##	Ozone	Solar.R	Wind	Temp	Month	Day
## 1	41	190	7.4	67	5	1
## 2	36	118	8.0	72	5	2
## 3	12	149	12.6	74	5	3
## 4	18	313	11.5	62	5	4
## 5	NA	NA	14.3	56	5	5
## 6	28	NA	14.9	66	5	6
## 7	23	299	8.6	65	5	7
## 8	19	99	13.8	59	5	8
## 9	8	19	20.1	61	5	9
## 10	NA	194	8.6	69	5	10

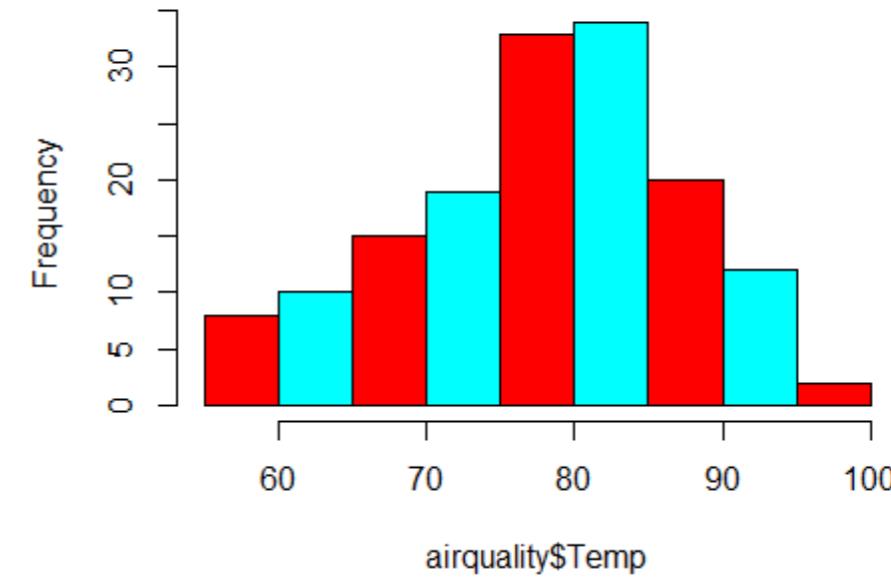
Month 5



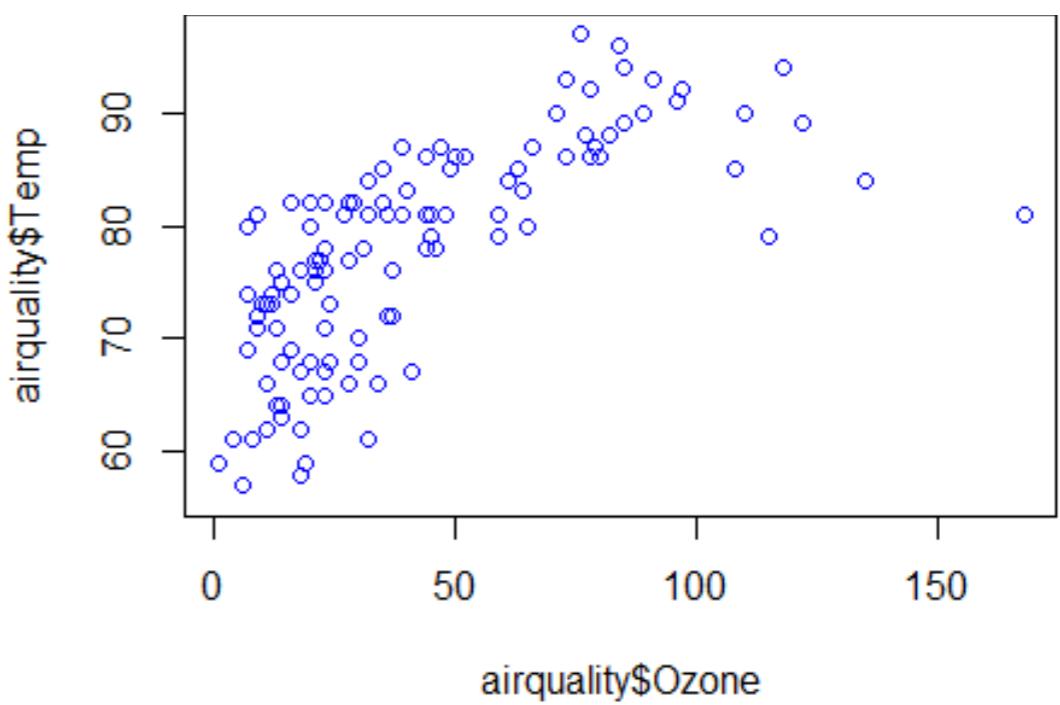
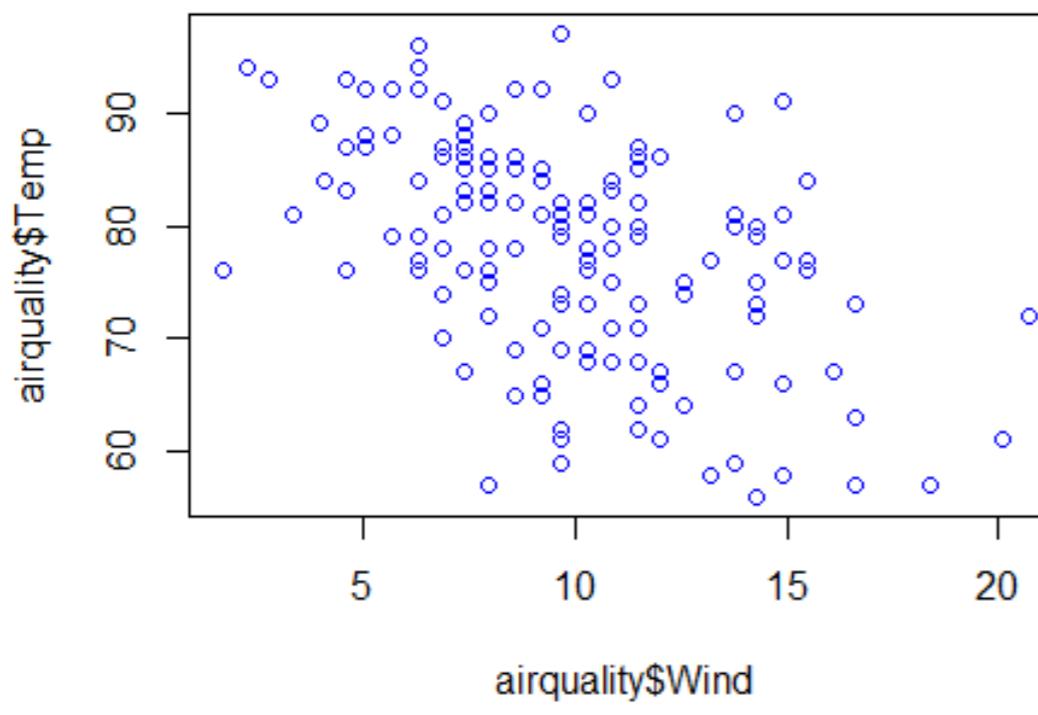
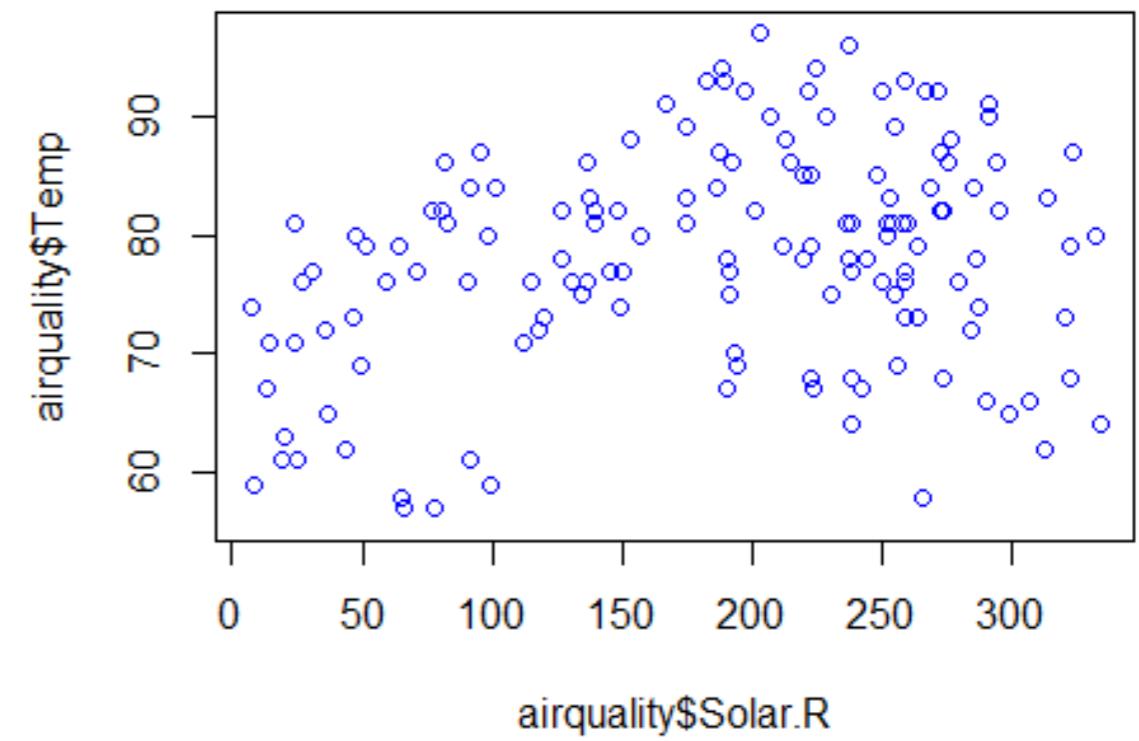
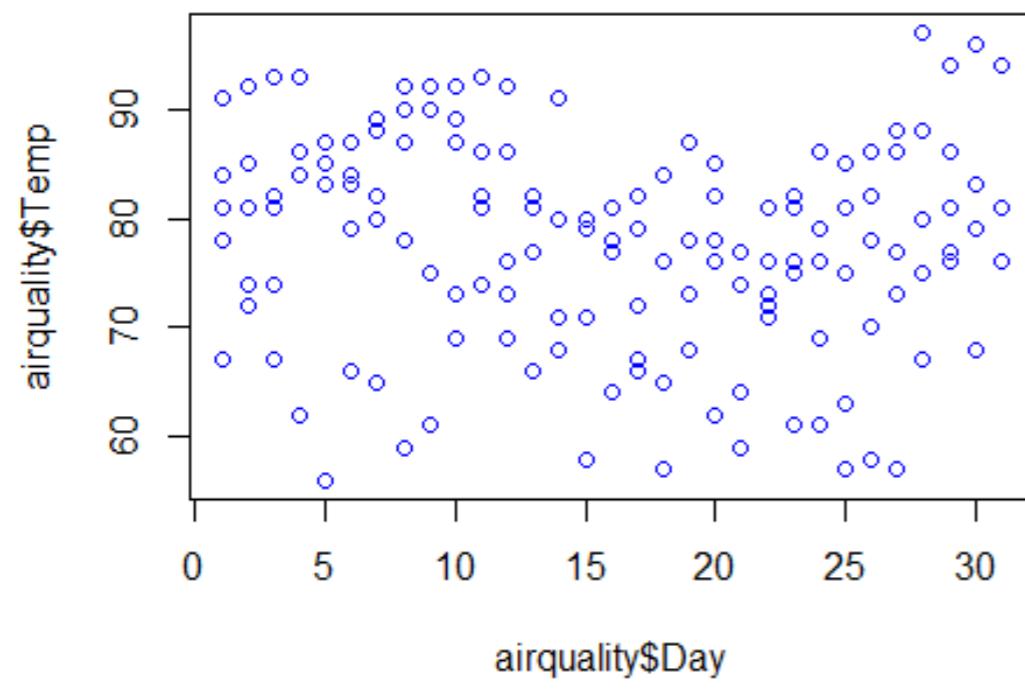
Month 6



Histogram of airquality\$Temp



Airquality data



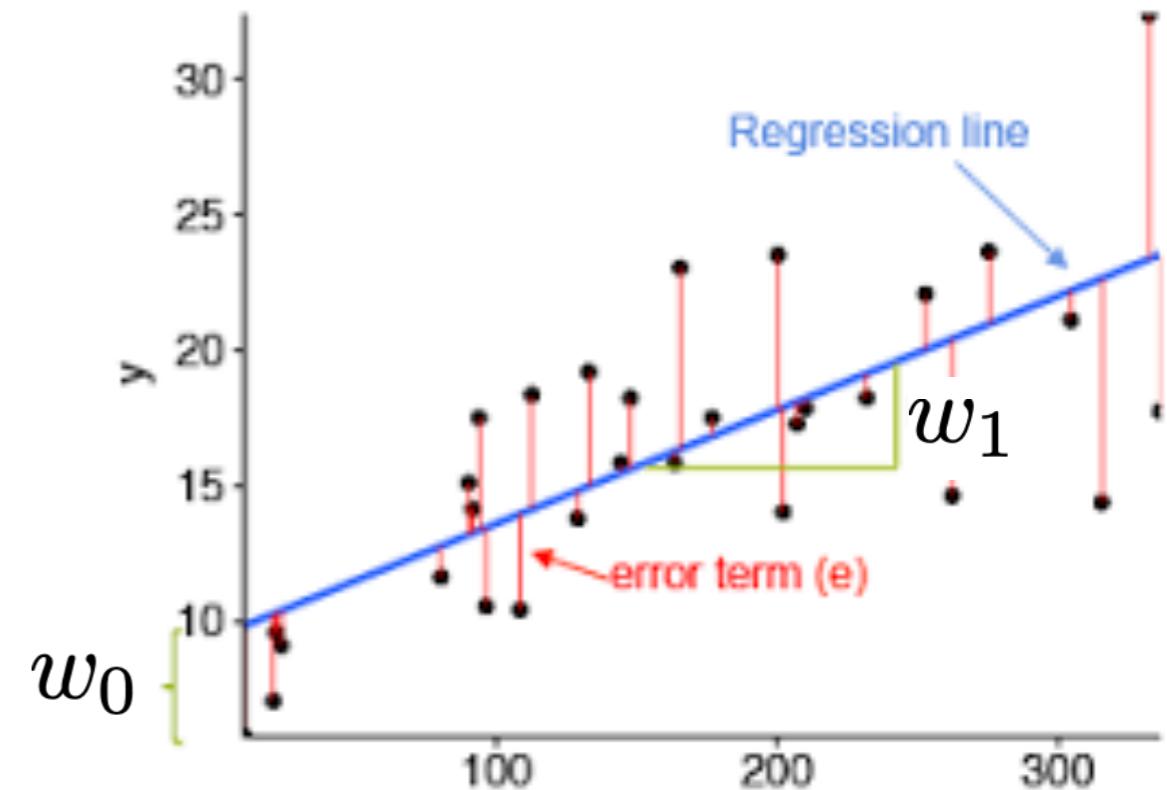
Linear regression

- $\text{Temp} = w_1 \cdot \text{Solar.R} + w_2 \cdot \text{Ozone} + w_3 \cdot \text{Wind} + \text{error}$.
- Temperature of house depends on ozone, wind and solar radiations
- linear regression helps to discover relation between dependent and independent variables

Linear regression

Regression Output is real and scalar, $y \in \mathbb{R}$

- Learn a function which maps input to output $f : X \rightarrow Y$
- Learn the function which passes through as many points as possible



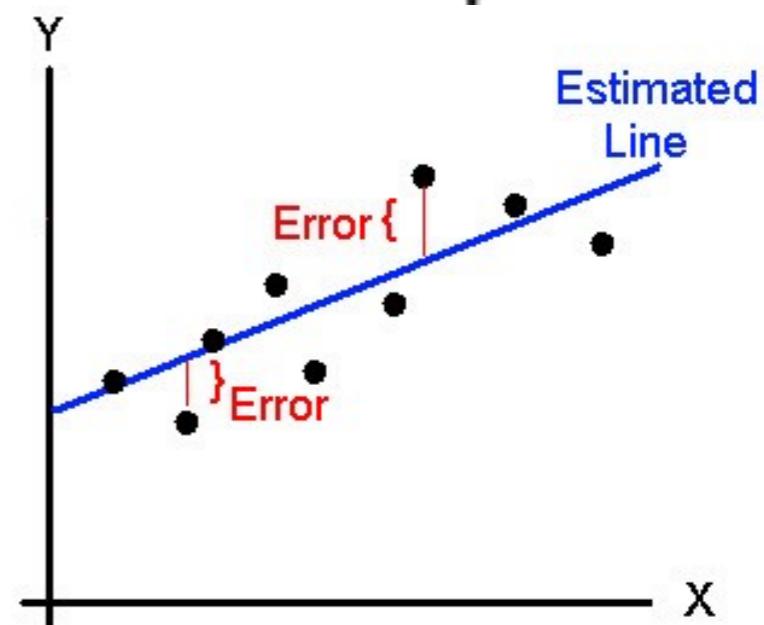
$$\hat{Y}_i = w_0 + w_1 X_i$$

Estimated (or predicted) Y value for observation i

Estimate of the regression intercept

Estimate of the regression slope

Value of X for observation i



Linear Regression

- Minimize the least squares error

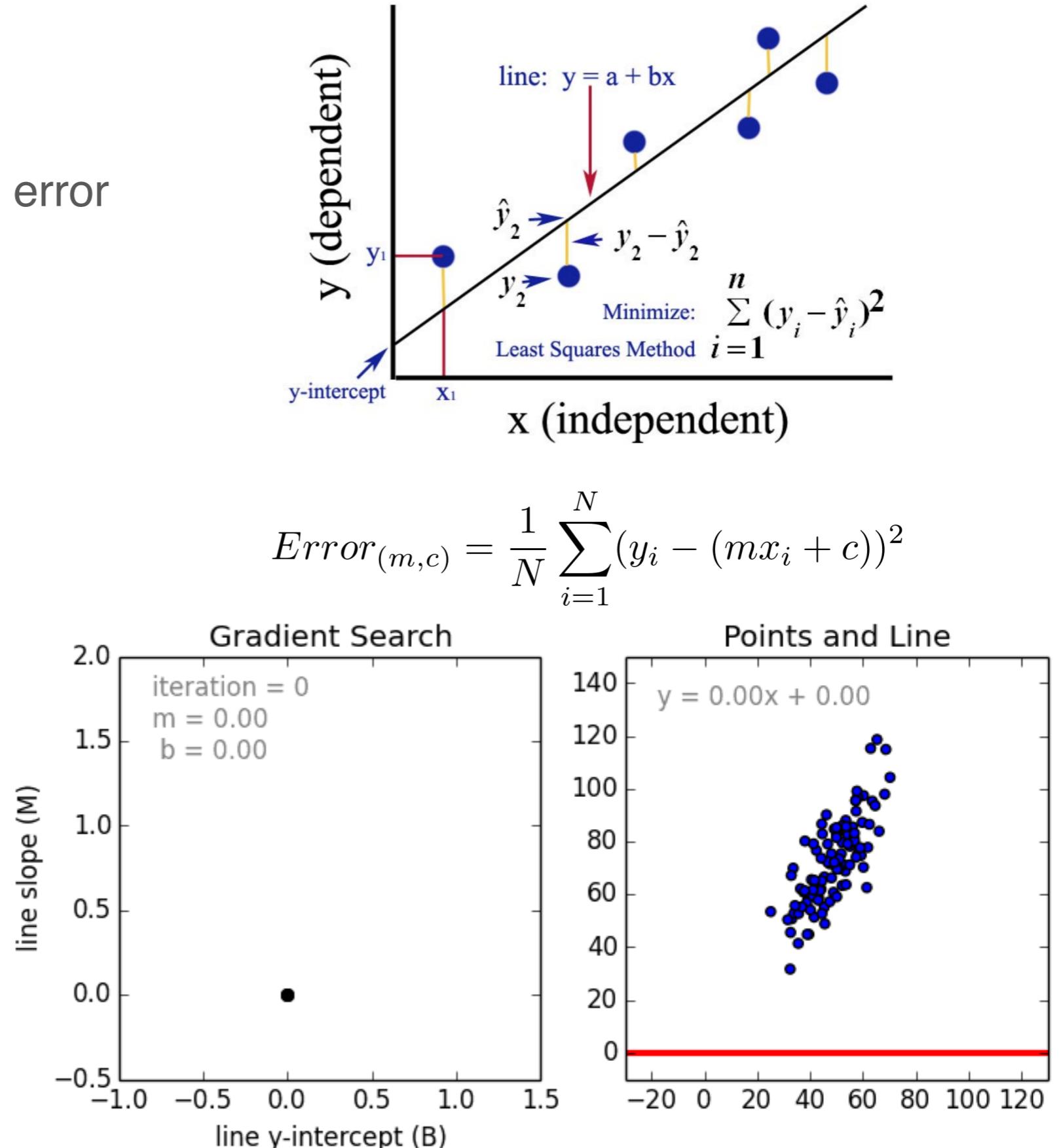
$$E(w) = \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$\begin{aligned} &= \frac{1}{2} \sum_{i=1}^N (y_i - X_i^\top w)^2 \\ &\quad \frac{1}{2} \|y - X^\top w\|^2 \end{aligned}$$

$$\nabla E(w) = Xy - XX^\top w = 0$$

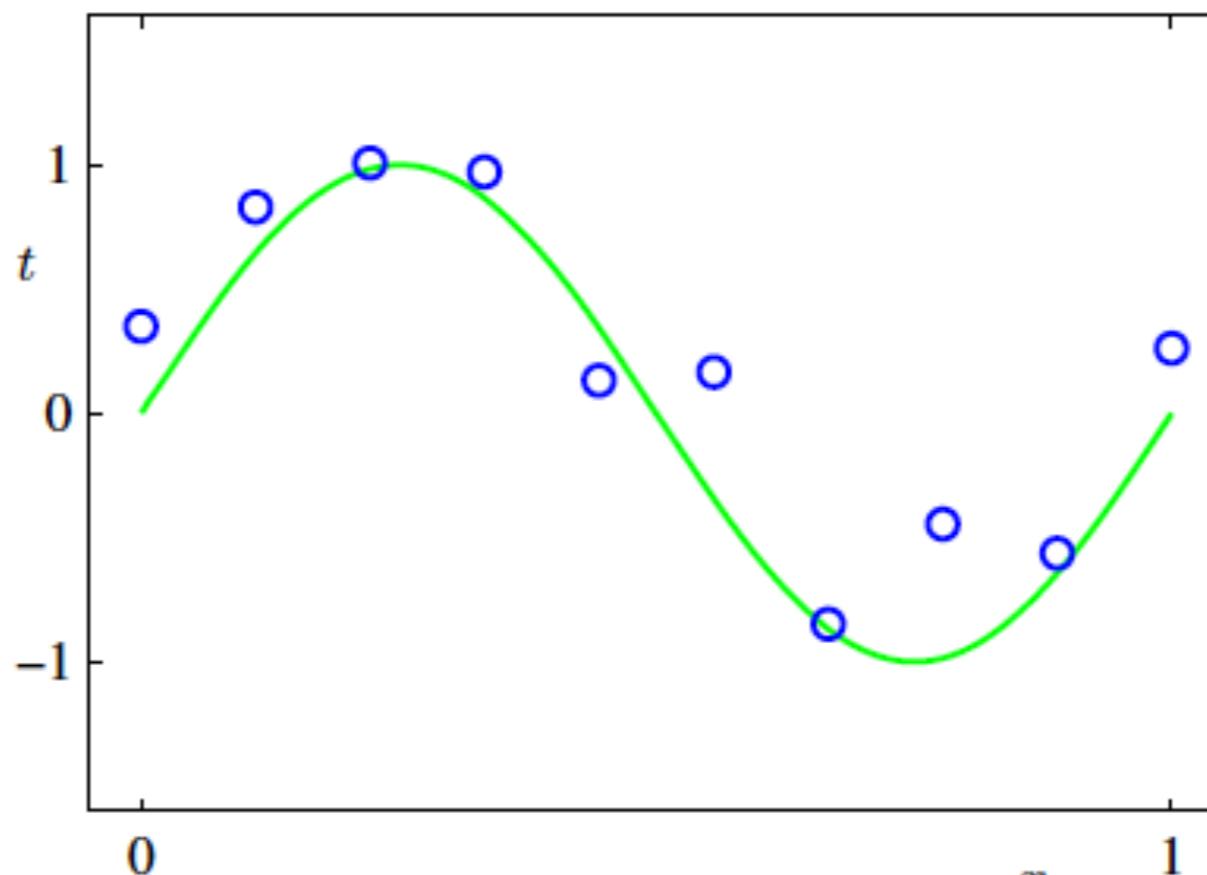
$$w_{ML} = (XX^\top)^{-1} Xy$$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n$$



Regression - curve fitting

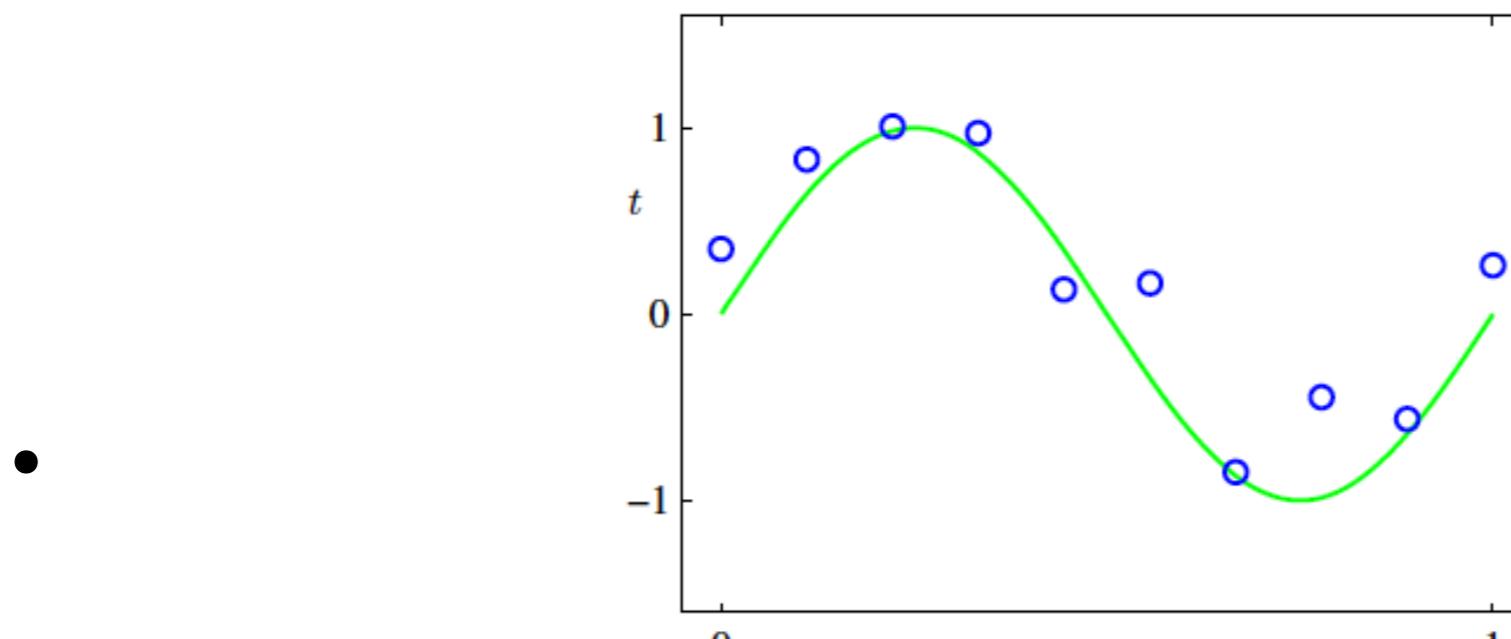
- Remember high school maths !
- Real-valued target variable t .
- Training set comprising N observations, features (input) $x \equiv (x_1, \dots, x_N)^T$, and target (output) $t \equiv (t_1, \dots, t_N)^T$.



Regression - curve fitting

- M is the order of the polynomial, $y(x,w)$ is a nonlinear function of x, it is a linear function of the coefficients w.
- Functions, such as the polynomial, which are linear in the unknown parameters have important properties and are called **linear models**

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

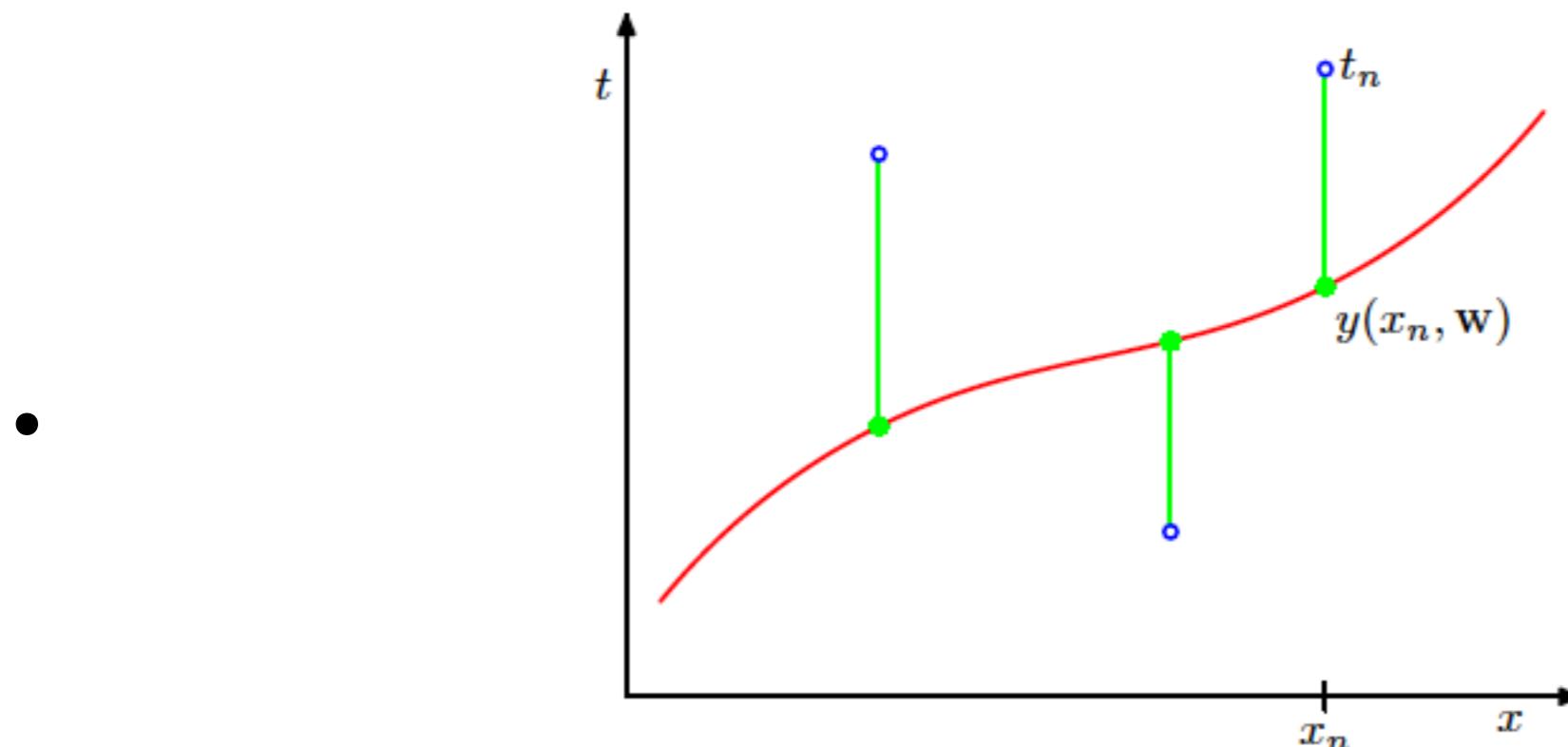


Regression - curve fitting

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

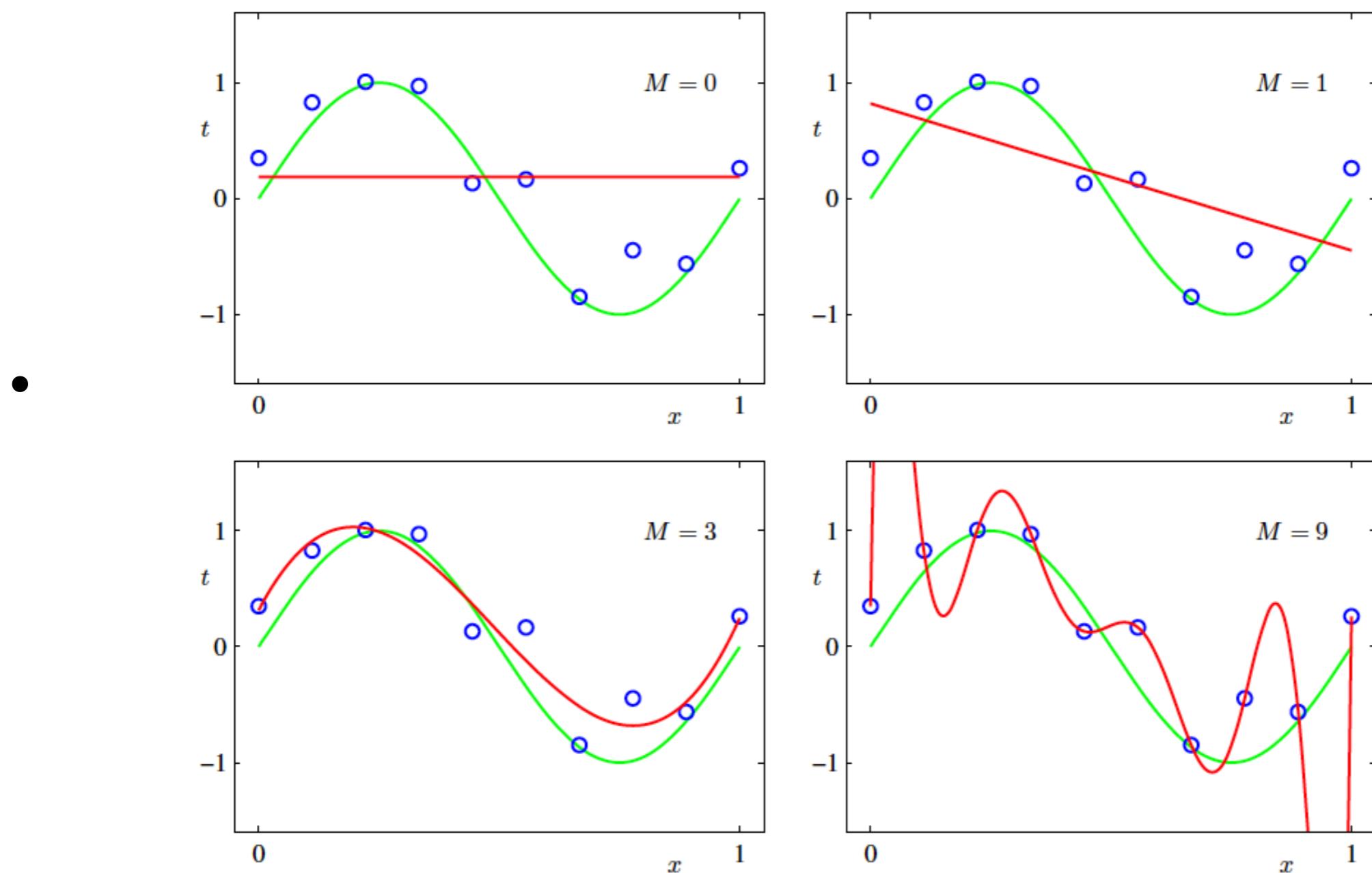
- Coefficients will be determined by fitting the polynomial to the training data. This can be done by minimizing an error function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$



Regression - curve fitting

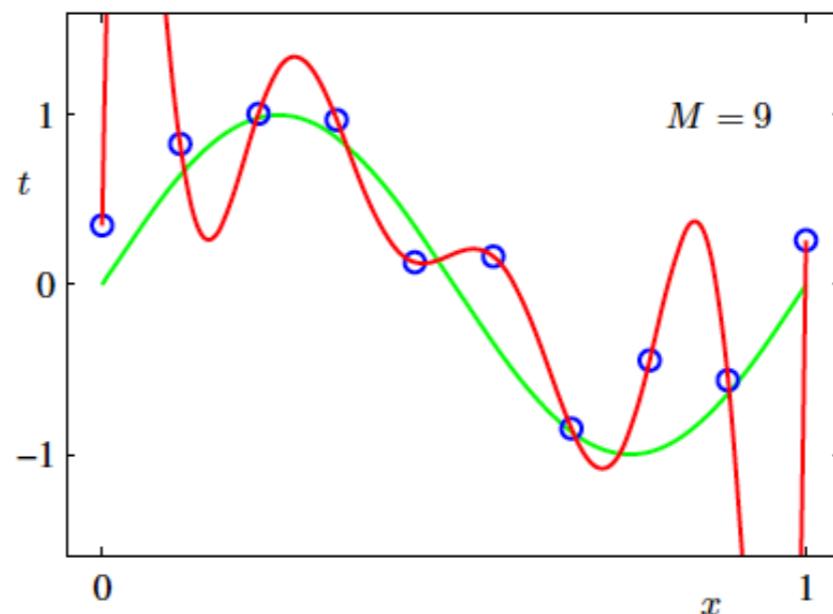
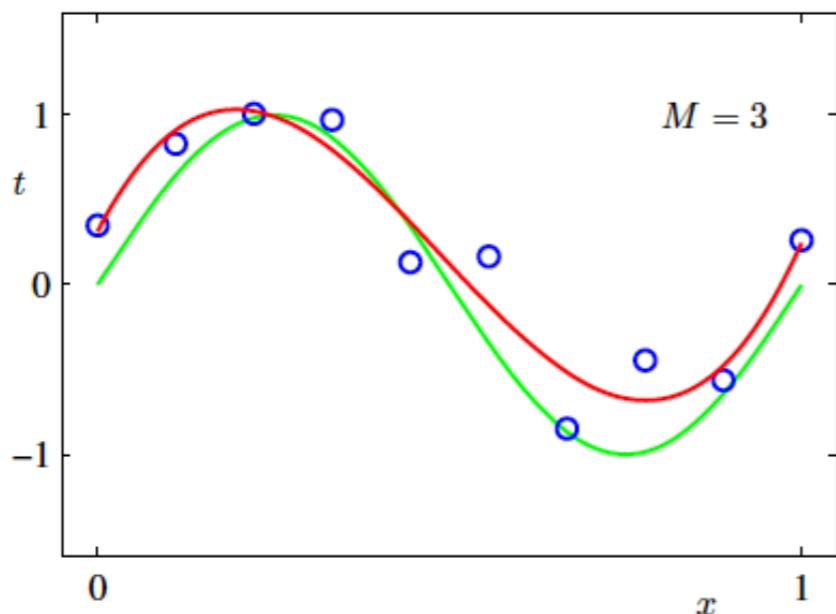
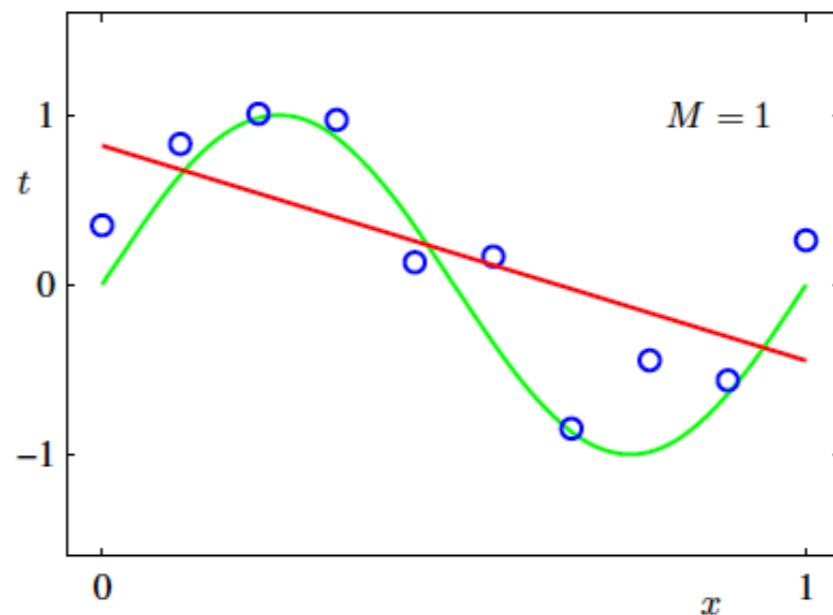
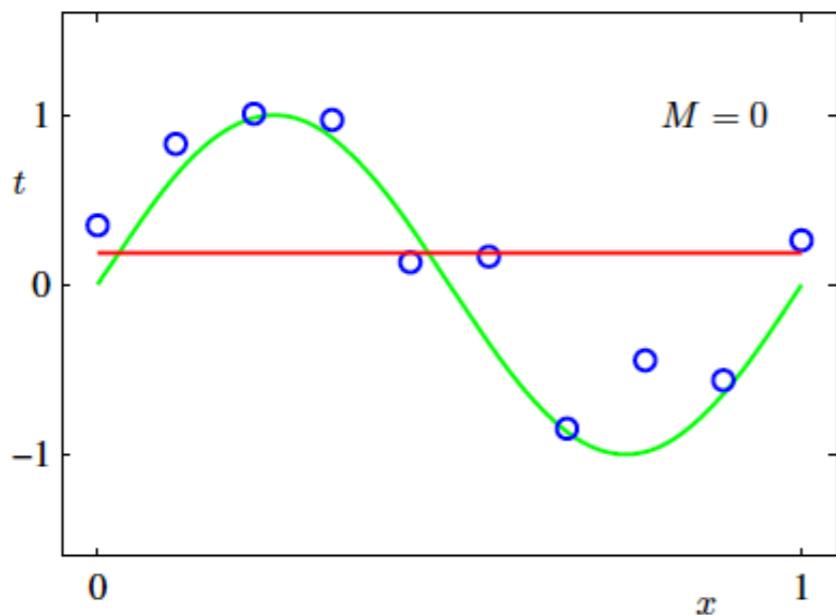
- Model selection (choosing M) : higher order polynomial ($M = 9$), provide excellent fit to the training data but gives a very poor representation of the function



Regression - curve fitting

- Model selection (choosing M) : high order polynomial ($M = 9$), provide excellent fit to the training data, a very poor representation of the function

Overfitting

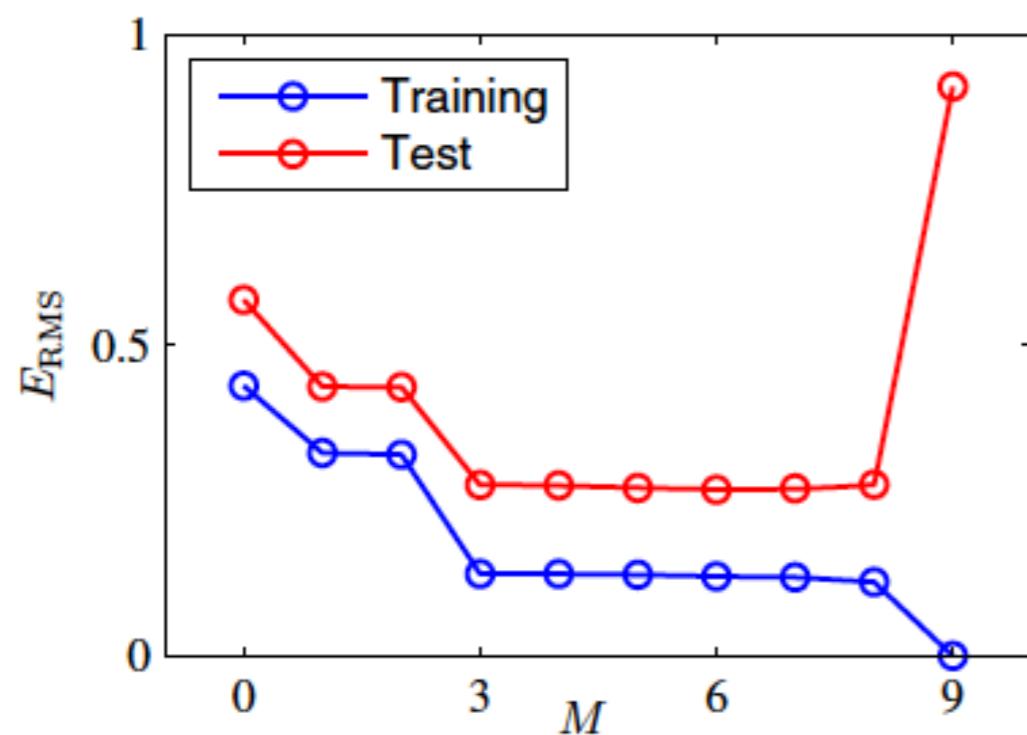


Regression - curve fitting

- Generalization performance measured using the root mean square error on test data

$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

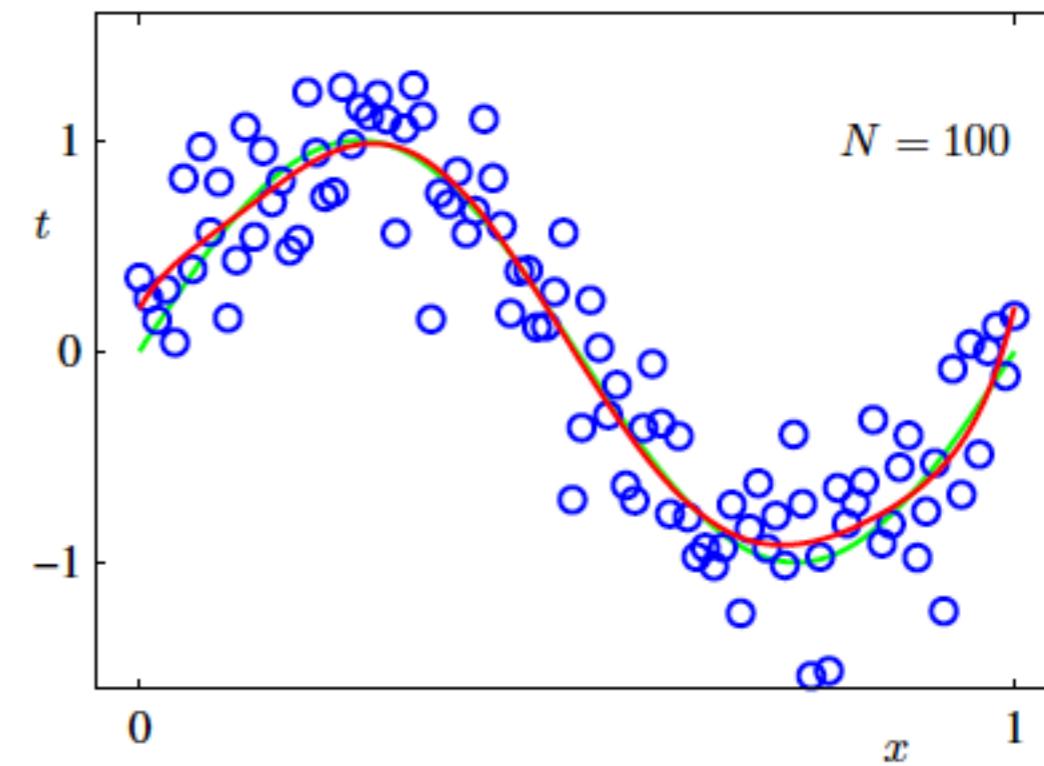
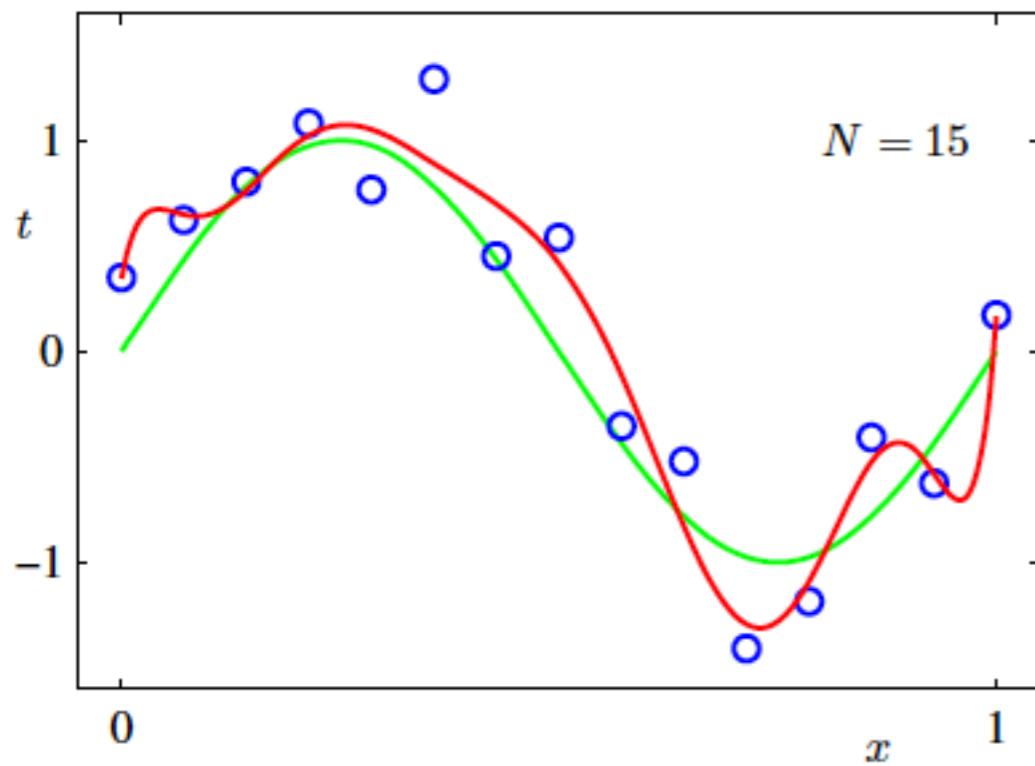


	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19		0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

- More flexible polynomials with larger values of M are becoming increasingly tuned to the random noise on the target values.

Regression - curve fitting

- Given model complexity, the over-fitting problem become less severe as the size of the data set increases.



Curve fitting - regularization

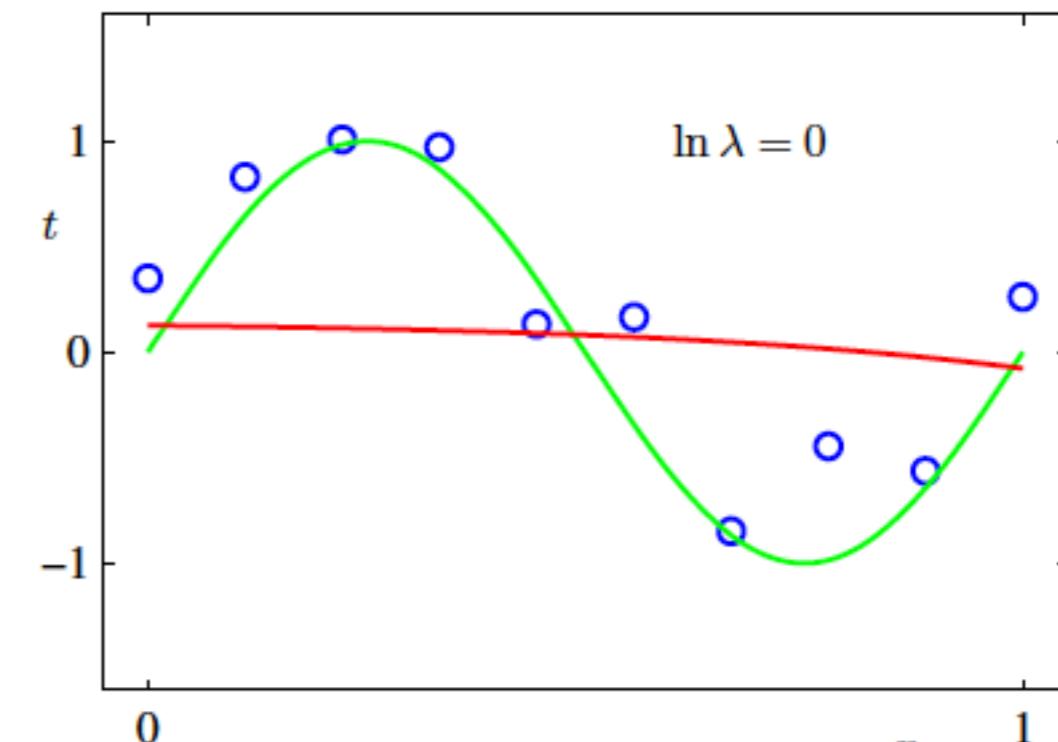
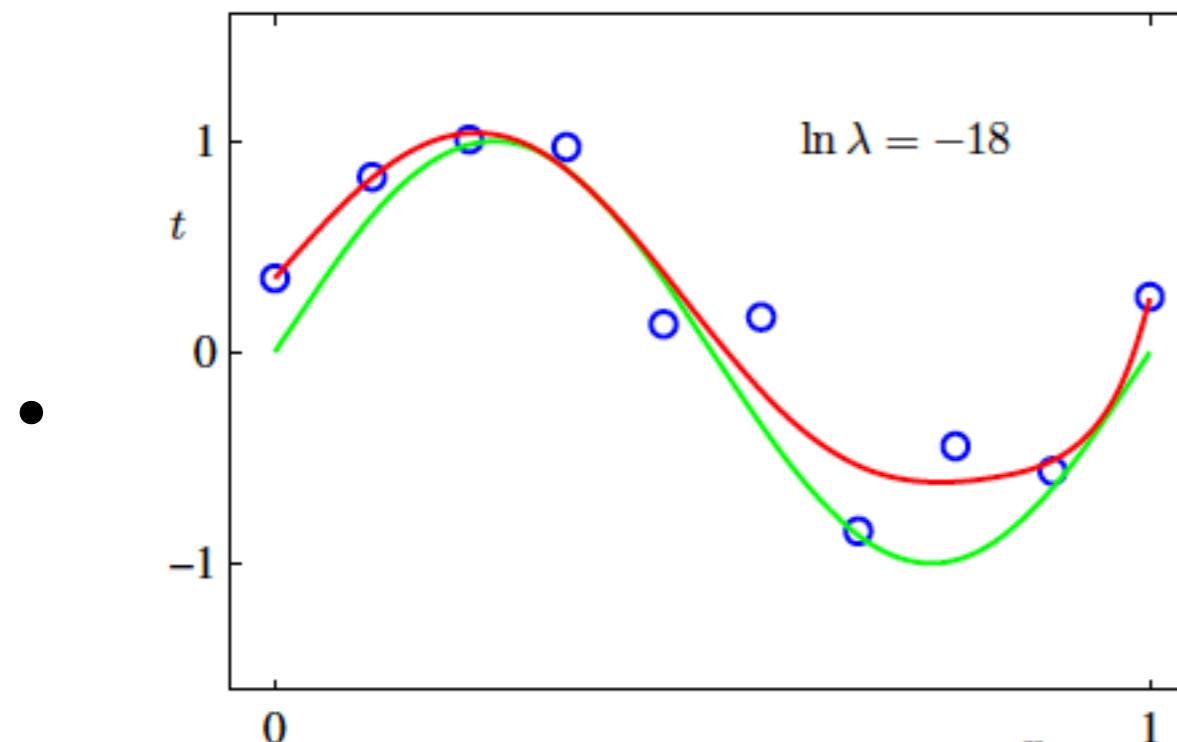
- Add a penalty term to the error function (1.2) in order to discourage the coefficients from reaching large values

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

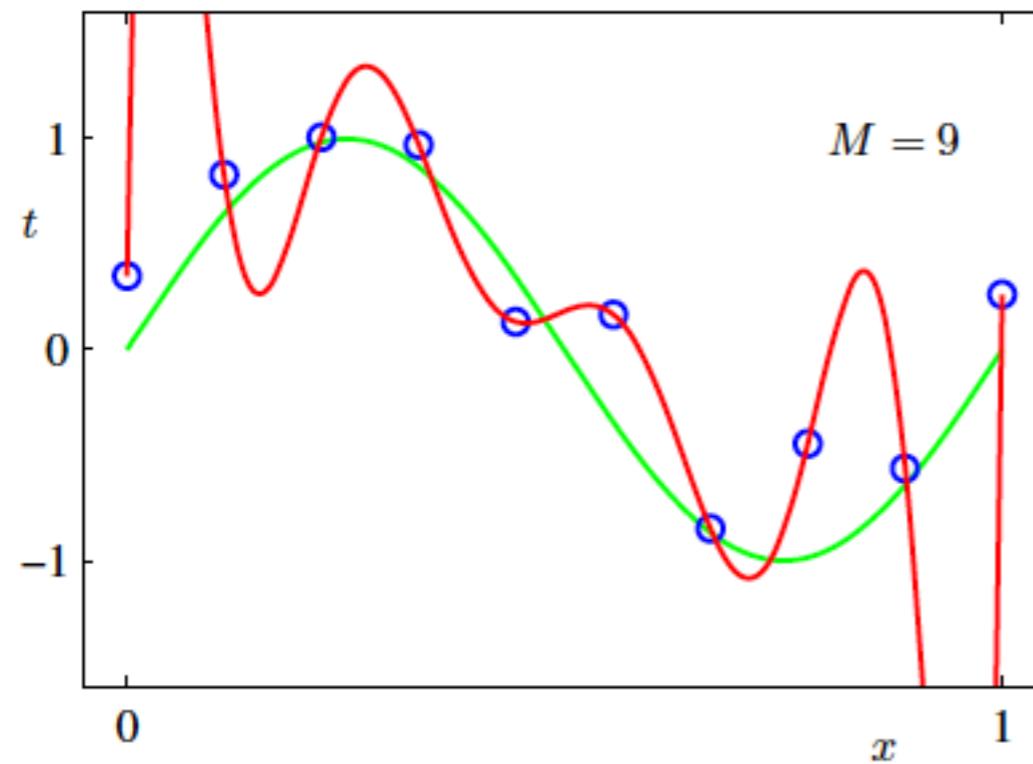
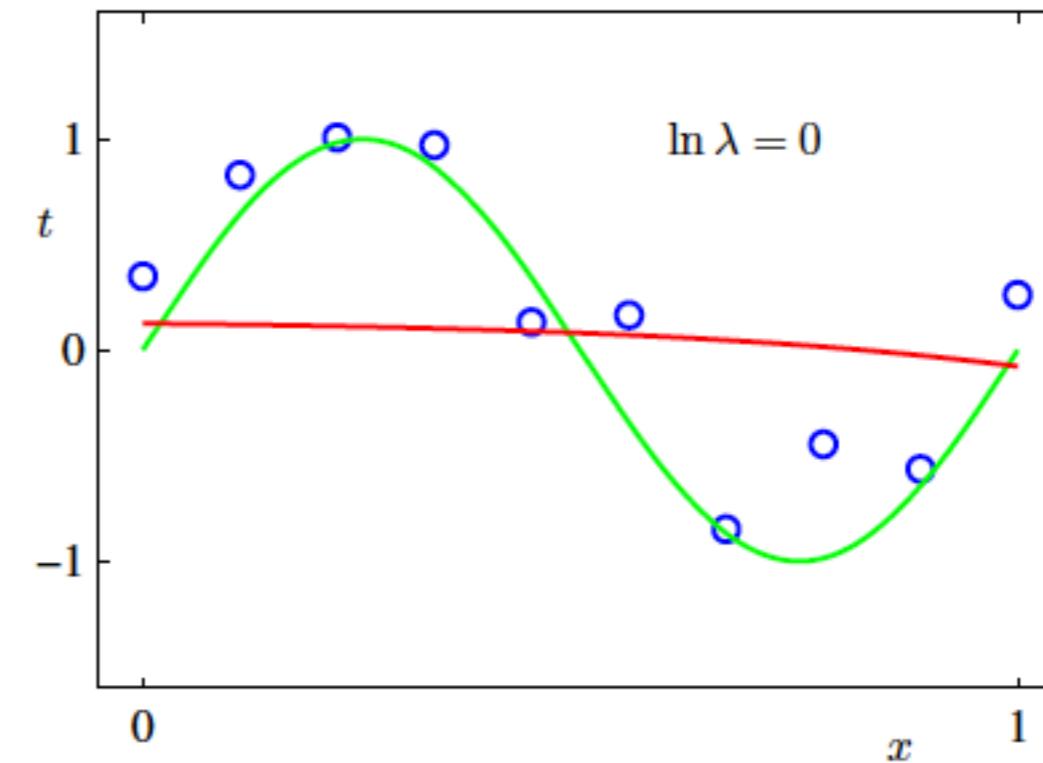
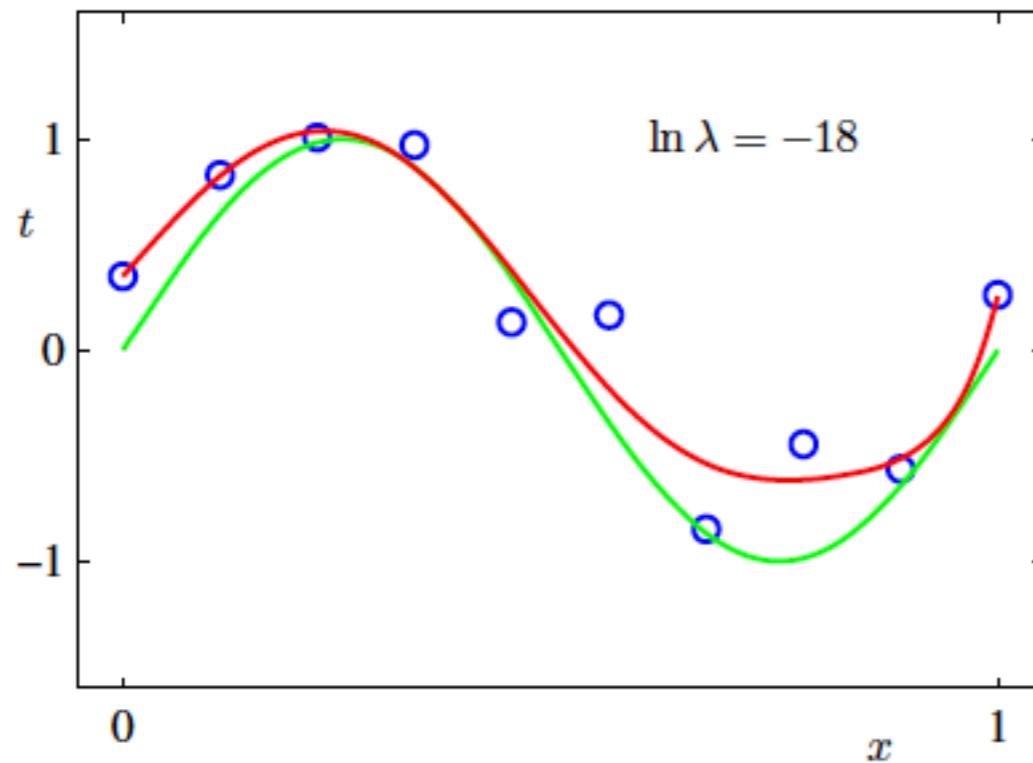
Regularization
constant

- Ridge regression : L2 norm

$$\|\mathbf{w}\|^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + \dots + w_M^2$$



Curve fitting - regularization



	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53	0.00
w_9^*	125201.43	72.68	0.01

Generalized linear models

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) \quad \mathbf{w} = (w_0, \dots, w_{M-1})^T \quad \phi = (\phi_0, \dots, \phi_{M-1})^T$$

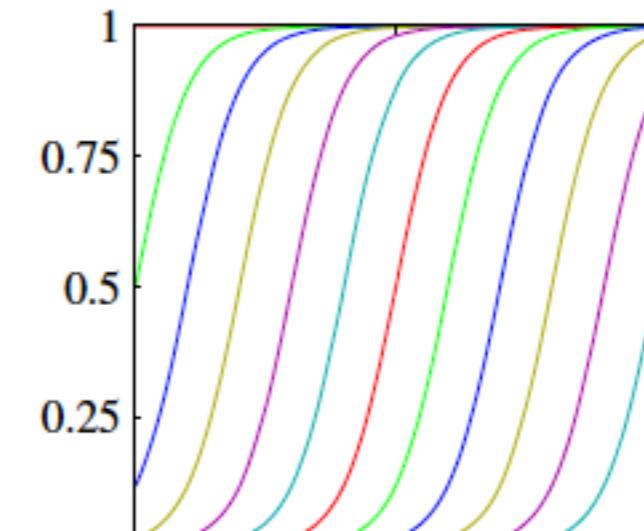
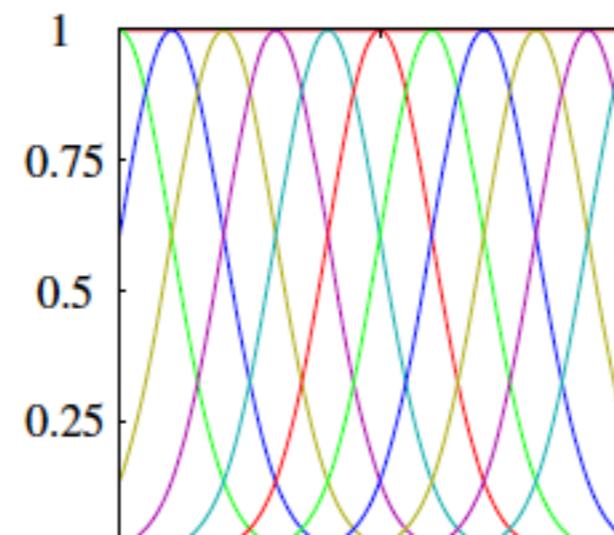
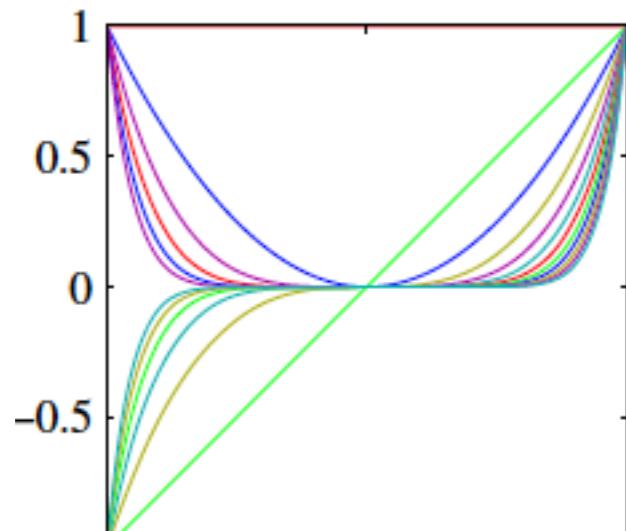
- $\phi_j(x)$ are known as basis functions, $\phi_j(x) = x^j$.

- Gaussian basis function

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

- Sigmoidal basis function

$$\phi_j(x) = \sigma \left(\frac{x - \mu_j}{s} \right) \quad \sigma(a) = \frac{1}{1 + \exp(-a)}.$$



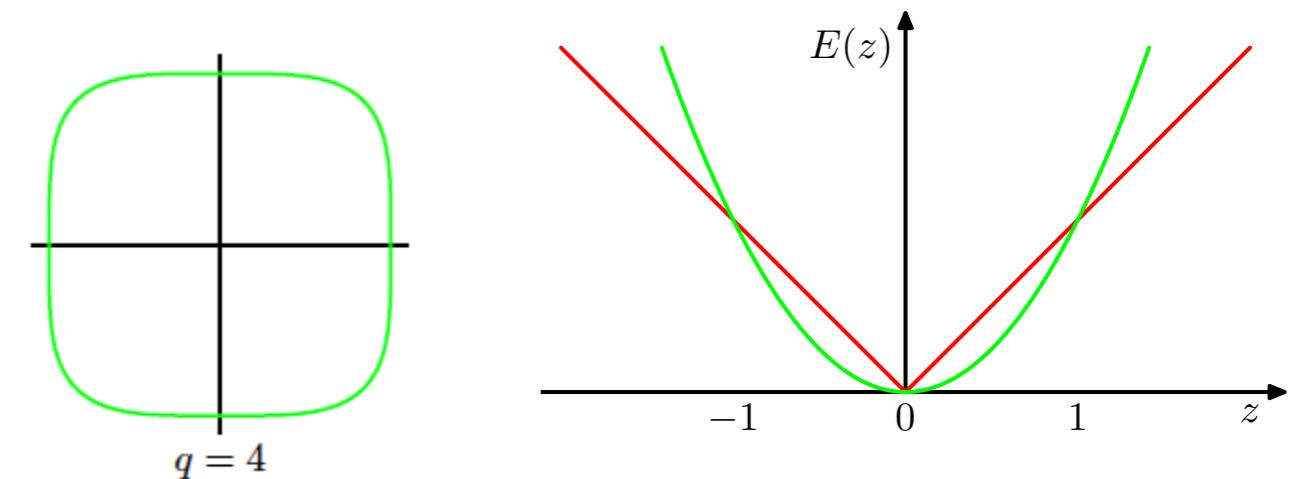
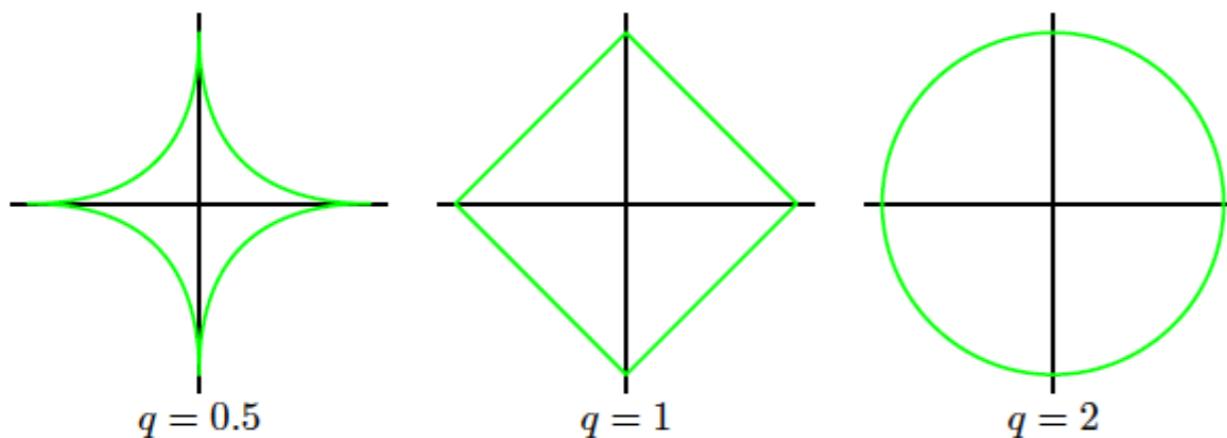
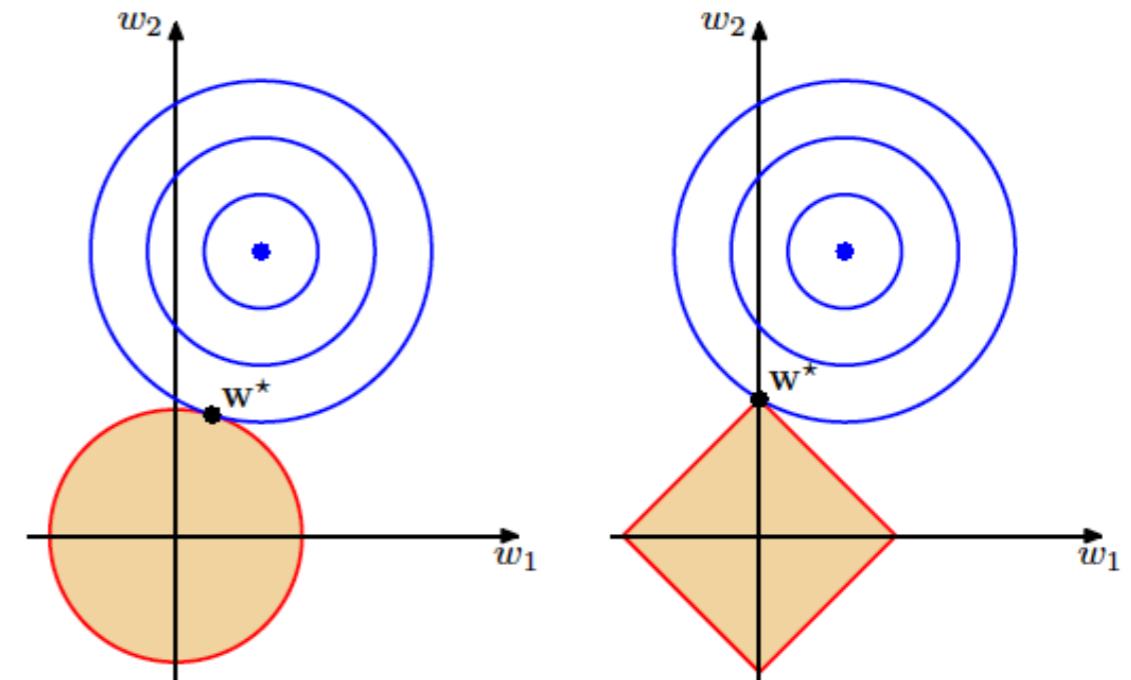
Regularized Least Squares

- Parameter shrinkage, weight decay
- **Ridge regression** $q=2$
- **Lasso regression** $q=1$, if λ is sufficiently large, some of the coefficients are driven to zero

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$$

$$\sum_{j=1}^M |w_j|^q \leq \eta$$

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$



Rain prediction : ML & MAP estimation

$$X = [1, 0, 0, 0, 1, 0, \dots]$$

Model ?

Likelihood : $p(X | \text{model})$

Learn model parameters :

Maximum Likelihood (ML) estimation

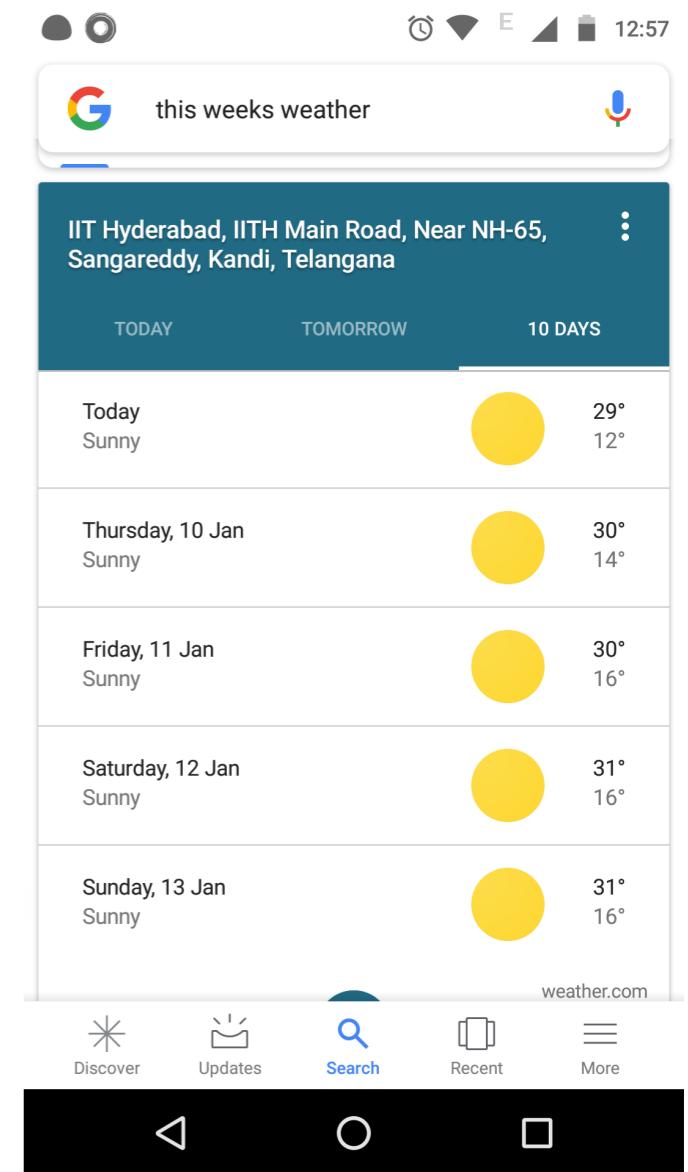
Maximum Aposteriori (MAP) estimation

$$\hat{\Theta}_{ML} = \operatorname{argmax}_{\Theta} \mathcal{L}$$

$$\mathcal{L} = \sum_{x_i \in \mathcal{X}} \log prob(x_i | \Theta)$$

$$\hat{\Theta}_{MAP} = \operatorname{argmax}_{\Theta} \prod_{x_i \in \mathcal{X}} prob(x_i | \Theta) \cdot prob(\Theta)$$

$$\hat{\Theta}_{MAP} = \operatorname{argmax}_{\Theta} \left(\sum_{x_i \in \mathcal{X}} \log prob(x_i | \Theta) + \log prob(\Theta) \right)$$

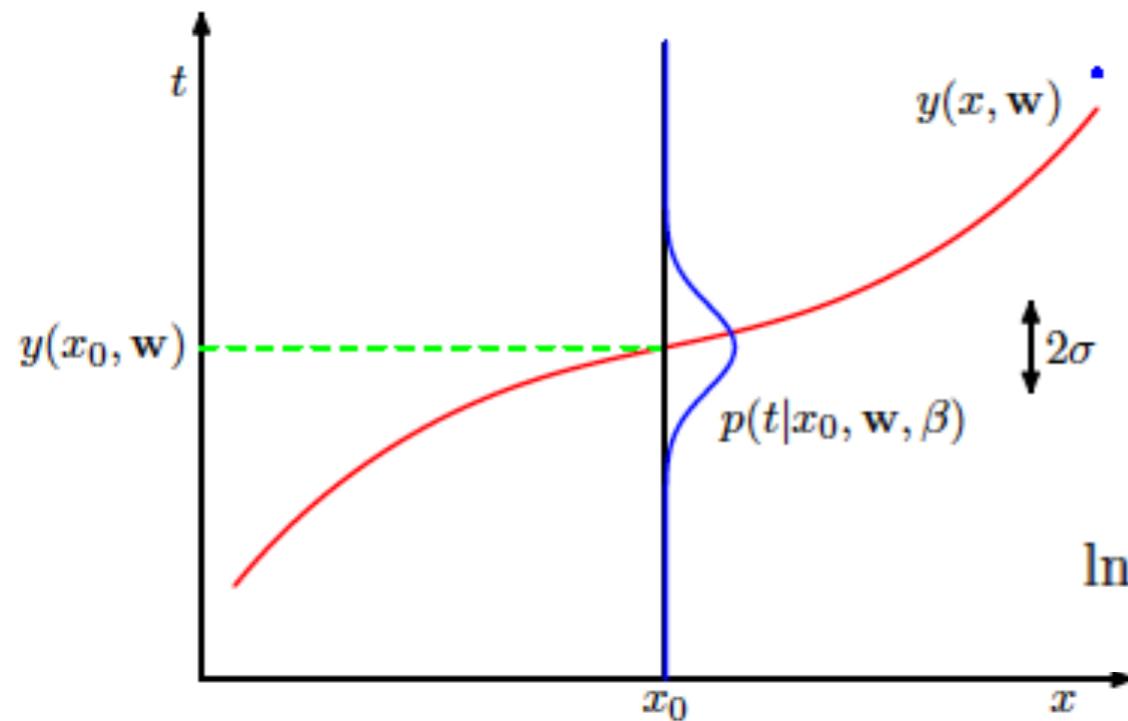


Curve fitting - A Probabilistic perspective

Least squares = maximum likelihood

- error function results from the maximum likelihood solution under an assumed Gaussian noise model

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1}) \quad t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2).$$



data points are i.i.d

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|y(x_n, \mathbf{w}), \beta^{-1}).$$

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi).$$

Curve fitting - A Probabilistic perspective

Least squares = maximum likelihood

- error function results from the maximum likelihood solution under an assumed Gaussian noise model

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi).$$

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2.$$

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T.$$

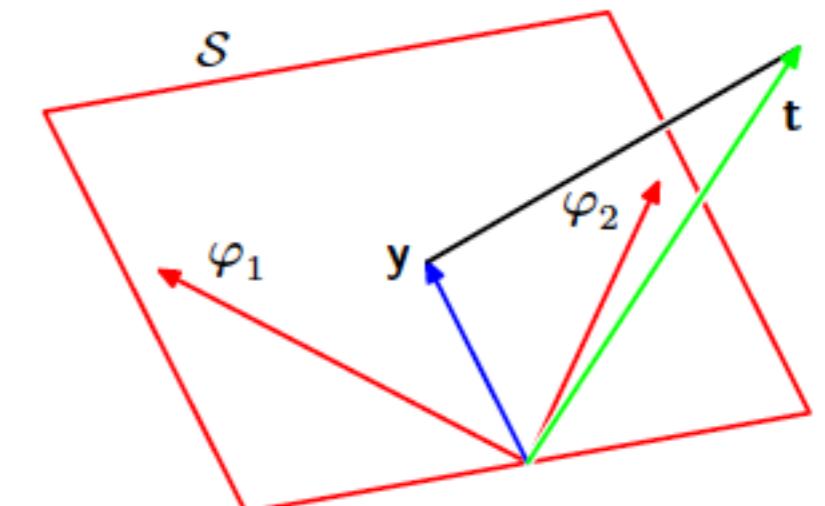
$$0 = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left(\sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right)$$

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

$$\Phi^\dagger \equiv (\Phi^T \Phi)^{-1} \Phi^T$$

Moore-Penrose pseudo-inverse



Regularized least squares regression = Maximum Aposteriori Estimate

- Ridge Regression

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- Compute Maximum aposteriori (MAP) estimate

- Prior over parameters

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

- Posterior

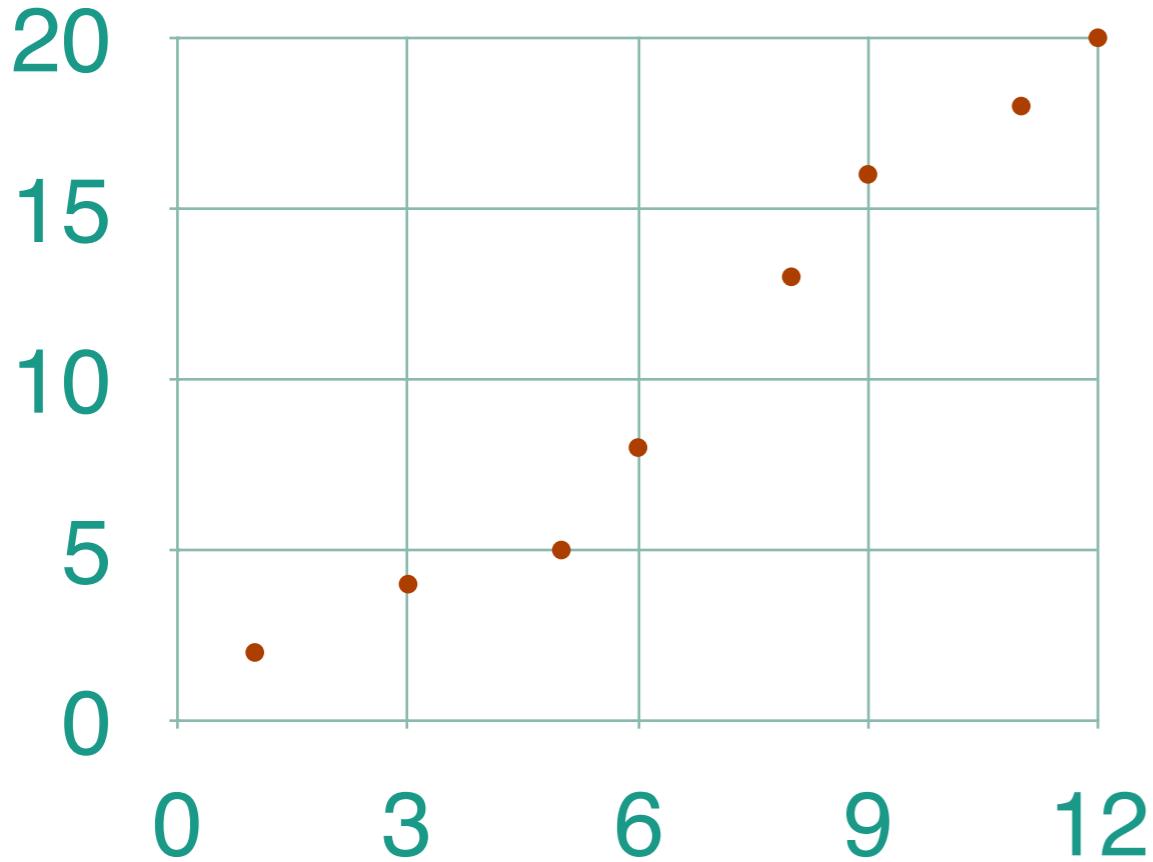
$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha).$$

- MAP estimate

$$\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}.$$

$$\mathbf{w} = (\lambda\mathbf{I} + \Phi^T\Phi)^{-1} \Phi^T \mathbf{t}. \quad \lambda = \alpha/\beta.$$

Bayesian Linear regression



- How to encode prior knowledge about the function e.g. slope
- Learn from limited data
- How to consider uncertainty
- Model selection
 - Which is the best function : linear, quadratic, cubic ?
 - regularisation constant
 - Cross-validation
 - Grid search is costly
 - only a fraction of data (validation data) is used for model selection
- Here comes Bayesian linear regression

Bayesian learning

Human Inductive Learning

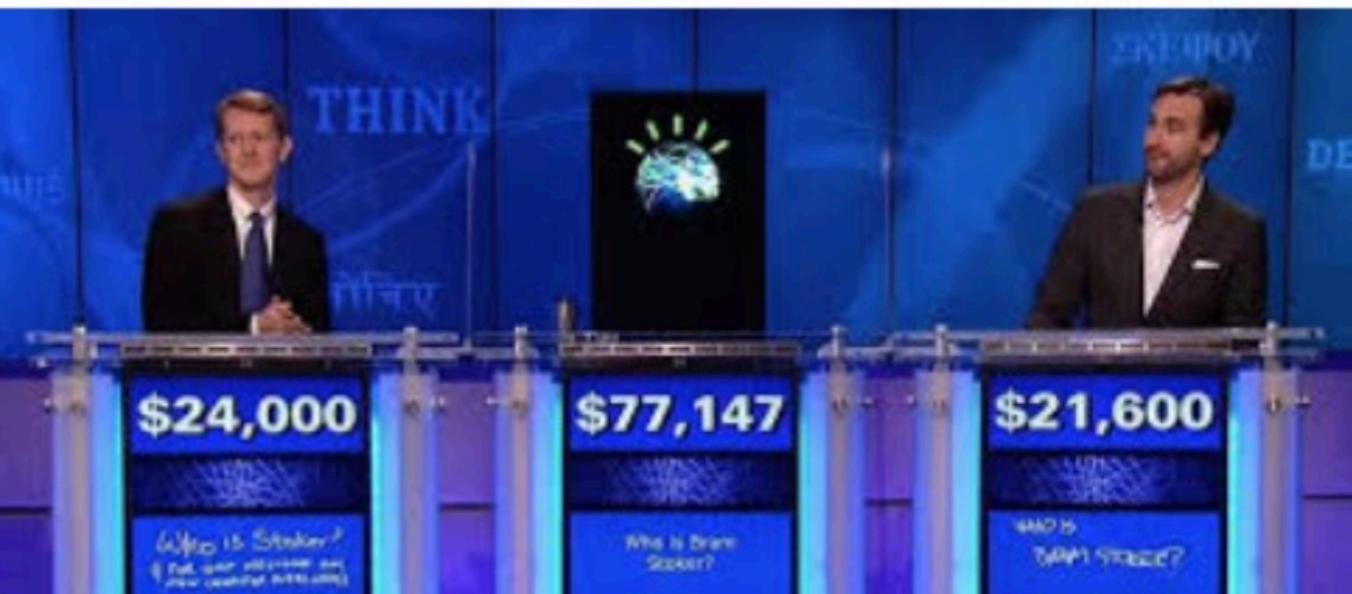


- how does the human mind go beyond the data of experience?
- how can it generalize well from the few noisy data ?

<https://gadgets.ndtv.com/science/news/lee-sedol-scores-surprise-victory-over-googles-alphago-in-game-4-813248>

Uncertainty in machine learning

- Safe AI : Uncertainty is important in Decision making
 - Important in high risk applications such as autonomous driving, disease diagnosis



Google Photos labeled black people 'gorillas'

Jessica Guynn, USA TODAY

Published 1:15 p.m. ET July 1, 2015 | Updated 2:10 p.m. ET July 1, 2015

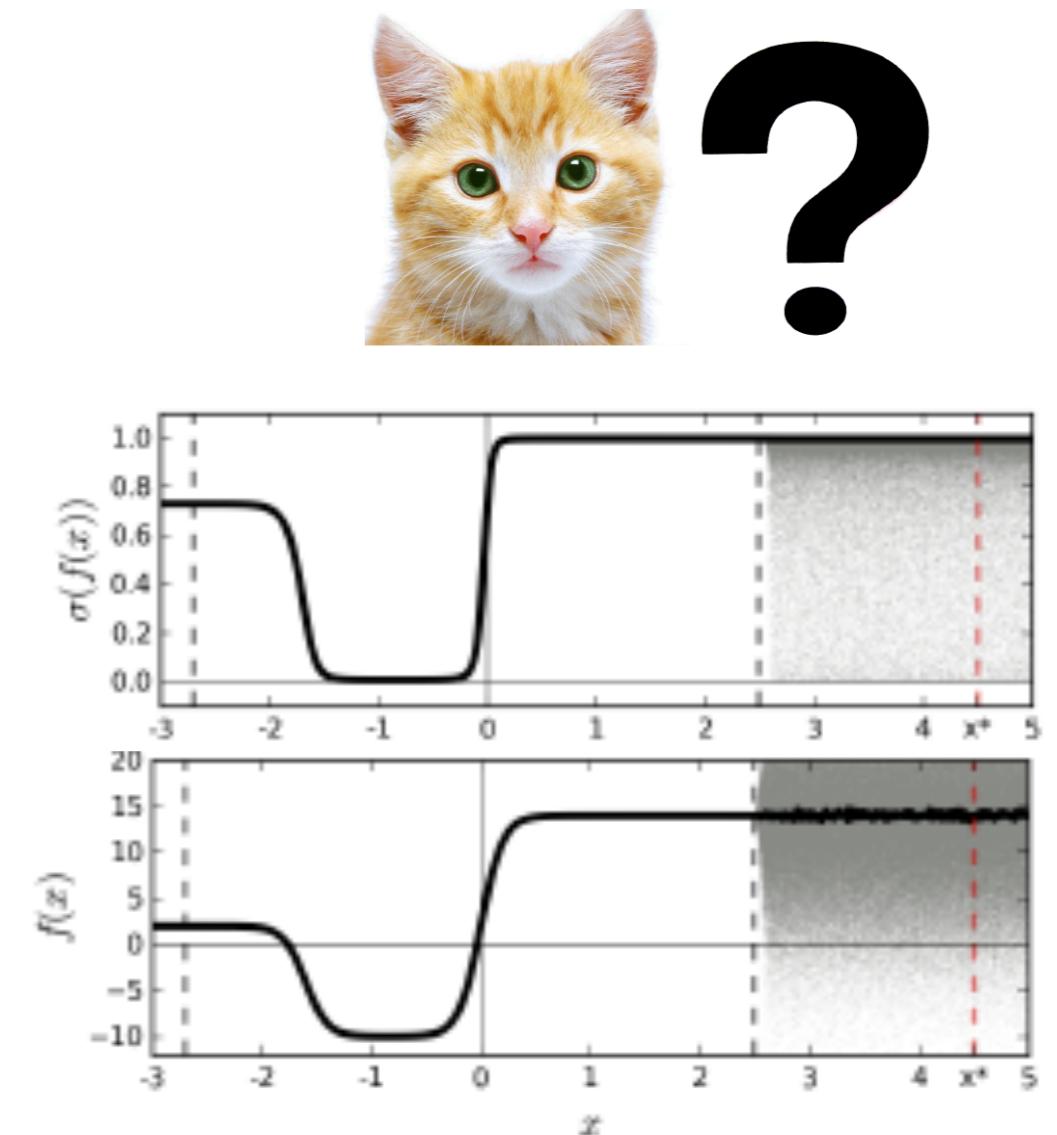
<https://www.topbots.com/10-major-fortune-500-brands-using-ibm-watson/>

<https://cleantechnica.com/2016/07/02/tesla-model-s-autopilot-crash-gets-bit-scary-negligent/>



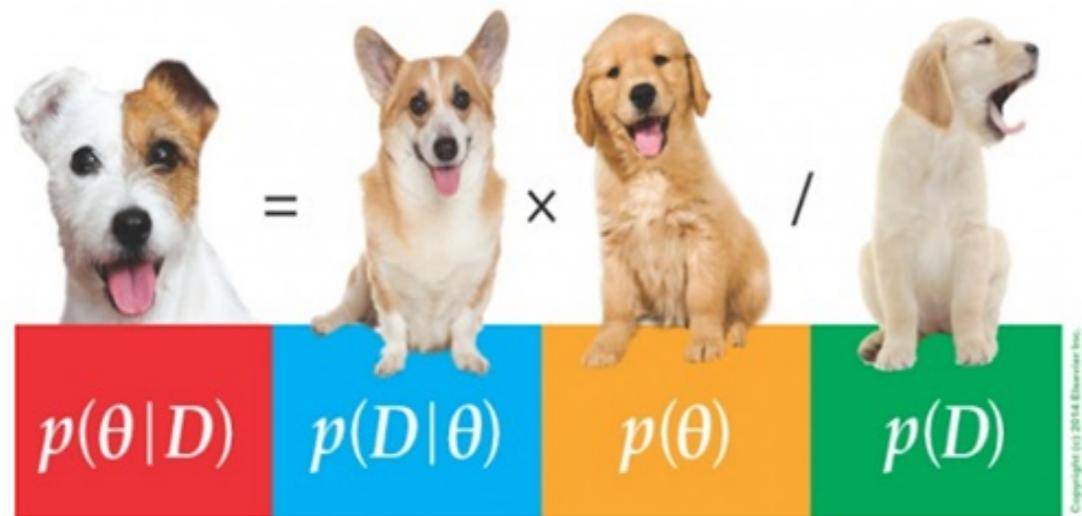
Uncertainty in machine learning

- Out of distribution data : Train the model on dogs breeds and provide picture of cat
- Uncertainty provides an idea on what the model actually know about the data

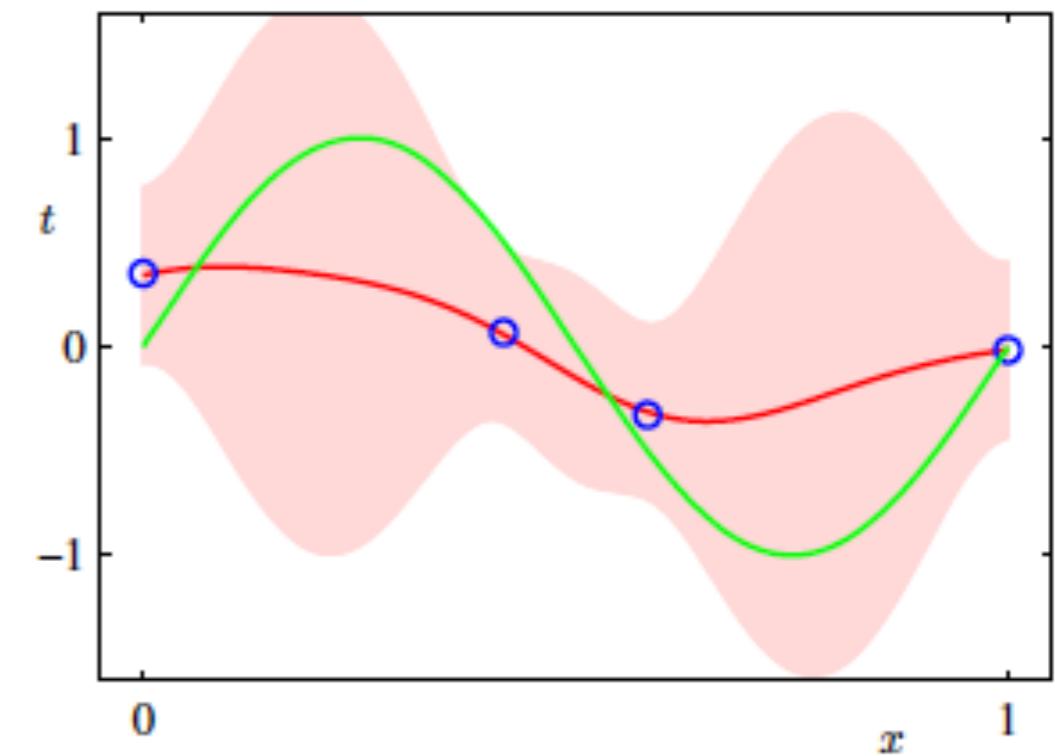


Bayesian learning

- Bayesian learning helps in better generalization capability and decision making by considering domain knowledge and by providing uncertainty estimates.



$$prob(\tilde{x}|\mathcal{X}) = \int_{\Theta} prob(\tilde{x}|\Theta) \cdot prob(\Theta|\mathcal{X}) d\Theta$$



- Bayesian linear regression, Gaussian processes, Latent Dirichlet allocation (LDA), Dirichlet Processes

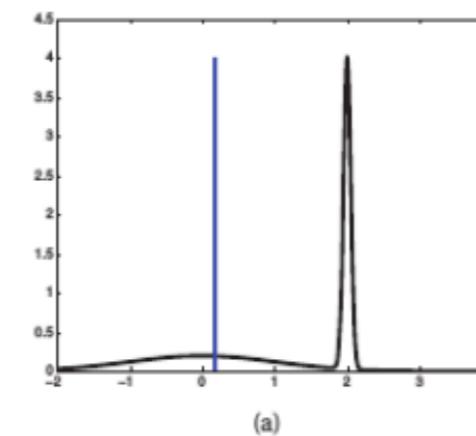
Rain prediction : Bayesian estimation

Both ML and MAP return only single and specific values for the parameter Θ .

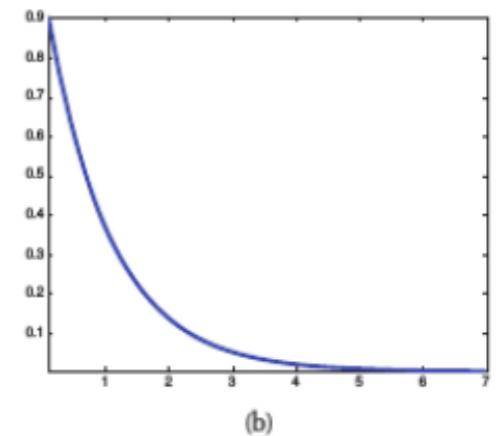
$$prob(\Theta|\mathcal{X}) = \frac{prob(\mathcal{X}|\Theta) \cdot prob(\Theta)}{prob(\mathcal{X})}$$

$$prob(\mathcal{X}) = \int_{\Theta} prob(\mathcal{X}|\Theta) \cdot prob(\Theta) d\Theta$$

$$prob(\tilde{x}|\mathcal{X}) = \int_{\Theta} prob(\tilde{x}|\Theta) \cdot prob(\Theta|\mathcal{X}) d\Theta$$



(a)



(b)

ML, MAP, Bayesian Estimate : Rain prediction

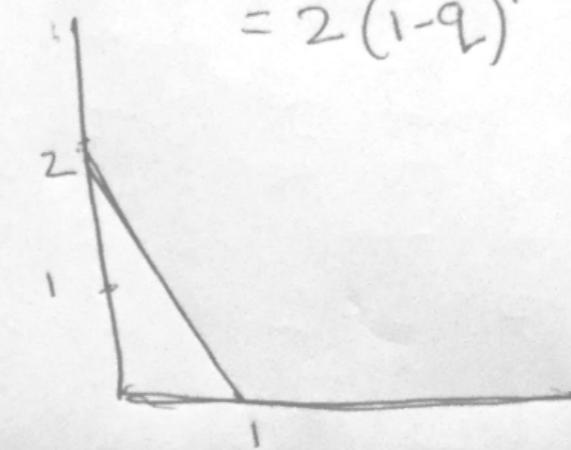
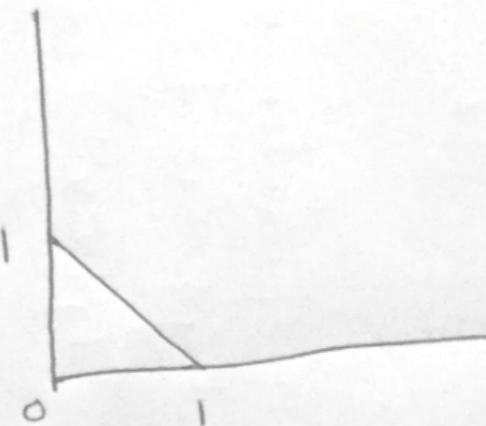
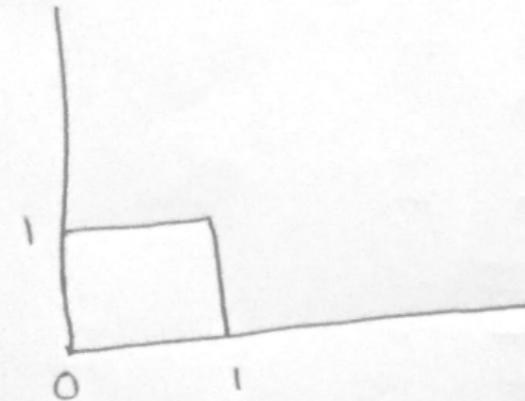
X : BINARY R.V.

$$X \sim \text{BERNOULLI}(q) \quad P(X|q) = q^x (1-q)^{1-x}$$

PRIOR $P(q) \sim \text{UNIF}(0,1)$

LIKELIHOOD $P(X=0|q) = 1-q$

POSTERIOR $P(q|X=0) = \frac{P(X=0|q) P(q)}{P(X)}$
 $= 2(1-q)^P(X)$



PREDICTION: $P(X^*=1|X) = \int p(X^*=1|q) P(q|X) dq$

Bayesian Machine learning : Bayesian linear regression

Learn $f : \mathcal{X} \rightarrow \mathcal{Y}$.

Bayesian approach : Allows to encode prior belief over functions

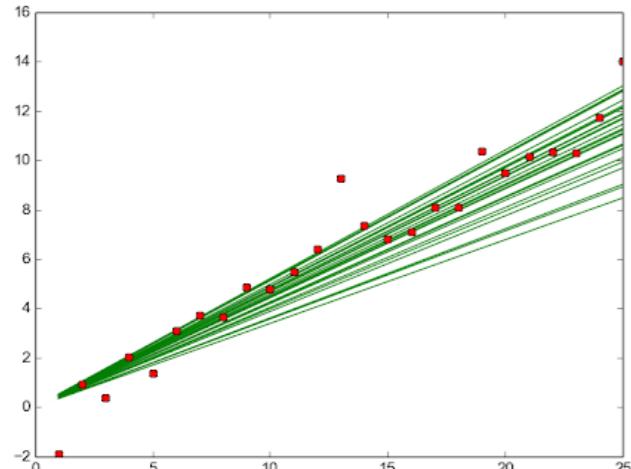
Parametric model : $f(x) = w \cdot x$

$$\begin{array}{ccc} p(w|D) & \propto & p(D|w) \\ \text{Posterior} & \propto & \text{Likelihood} \end{array} \quad p(w) \quad \text{Prior}$$

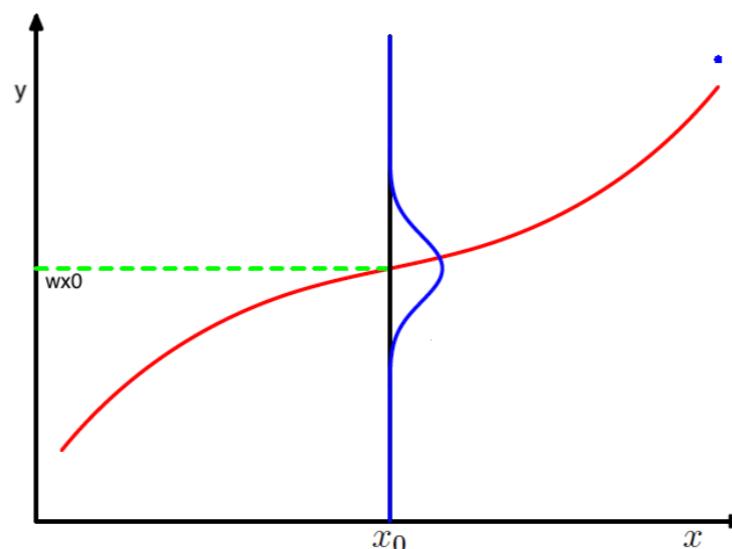
$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}, \quad y = f(\mathbf{x}) + \varepsilon,$$

$$\varepsilon \sim \mathcal{N}(0, \sigma_n^2).$$

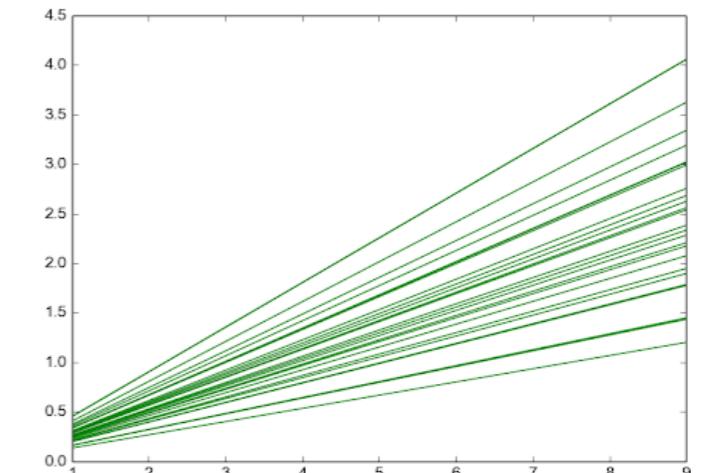
$$p(w|y, x) \propto p(y|w, x)p(w)$$



$$w \sim N(w; 0.45, 0.05)$$



$$p(y|w, x) = N(y; wx, \sigma^2)$$



$$w \sim N(w; 0.3, 0.1)$$

Bayesian linear regression

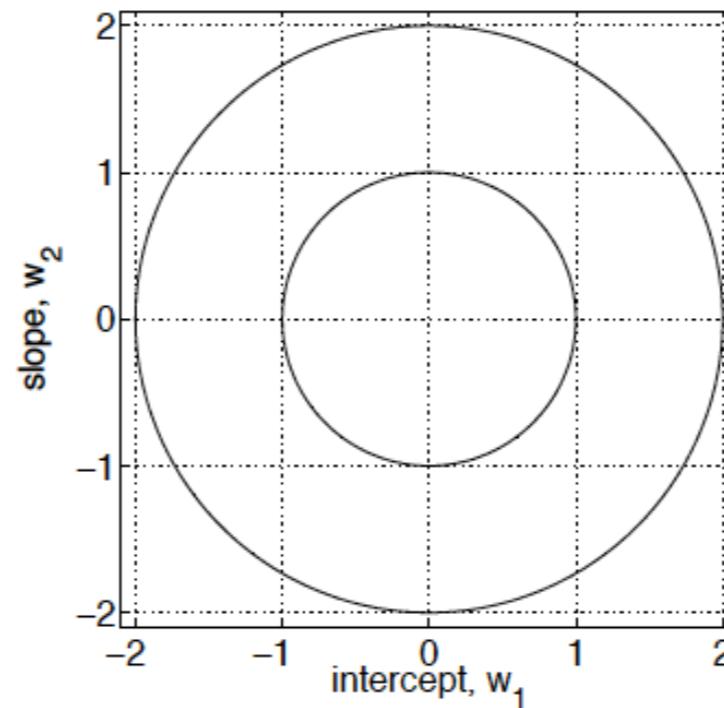
- Prior : $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)$.
- Likelihood :
$$\begin{aligned} p(\mathbf{y}|X, \mathbf{w}) &= \prod_{i=1}^n p(y_i|\mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left(-\frac{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}{2\sigma_n^2}\right) \\ &= \frac{1}{(2\pi\sigma_n^2)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2} |\mathbf{y} - X^\top \mathbf{w}|^2\right) = \mathcal{N}(X^\top \mathbf{w}, \sigma_n^2 I), \end{aligned}$$
- Maximum likelihood estimate $w_{ML} = (X X^\top)^{-1} X \mathbf{y}$
- Maximum a-posteriori estimate $\bar{\mathbf{w}} = \sigma_n^{-2} (\sigma_n^{-2} X X^\top + \Sigma_p^{-1})^{-1} X \mathbf{y}$

Bayesian linear regression

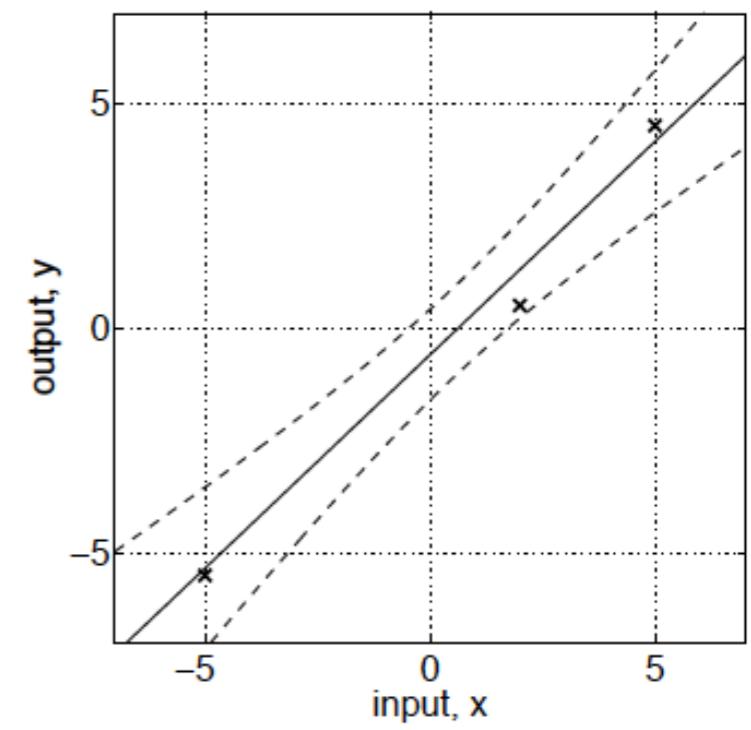
- Prior : $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)$.
- Likelihood :
$$\begin{aligned} p(\mathbf{y}|X, \mathbf{w}) &= \prod_{i=1}^n p(y_i|\mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left(-\frac{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}{2\sigma_n^2}\right) \\ &= \frac{1}{(2\pi\sigma_n^2)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2} |\mathbf{y} - X^\top \mathbf{w}|^2\right) = \mathcal{N}(X^\top \mathbf{w}, \sigma_n^2 I), \end{aligned}$$
- Posterior :
$$\begin{aligned} p(\mathbf{w}|X, \mathbf{y}) &\propto \exp\left(-\frac{1}{2\sigma_n^2} (\mathbf{y} - X^\top \mathbf{w})^\top (\mathbf{y} - X^\top \mathbf{w})\right) \exp\left(-\frac{1}{2} \mathbf{w}^\top \Sigma_p^{-1} \mathbf{w}\right) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{w} - \bar{\mathbf{w}})^\top \left(\frac{1}{\sigma_n^2} XX^\top + \Sigma_p^{-1}\right) (\mathbf{w} - \bar{\mathbf{w}})\right) \end{aligned}$$
$$p(\mathbf{w}|X, \mathbf{y}) \sim \mathcal{N}(\bar{\mathbf{w}} = \frac{1}{\sigma_n^2} A^{-1} X \mathbf{y}, A^{-1}), \quad A = \sigma_n^{-2} XX^\top + \Sigma_p^{-1}.$$
- Mean, MAP (prior mean 0) and ML estimates (prior variance large)

$$f(x) = w_1 + w_2 x$$

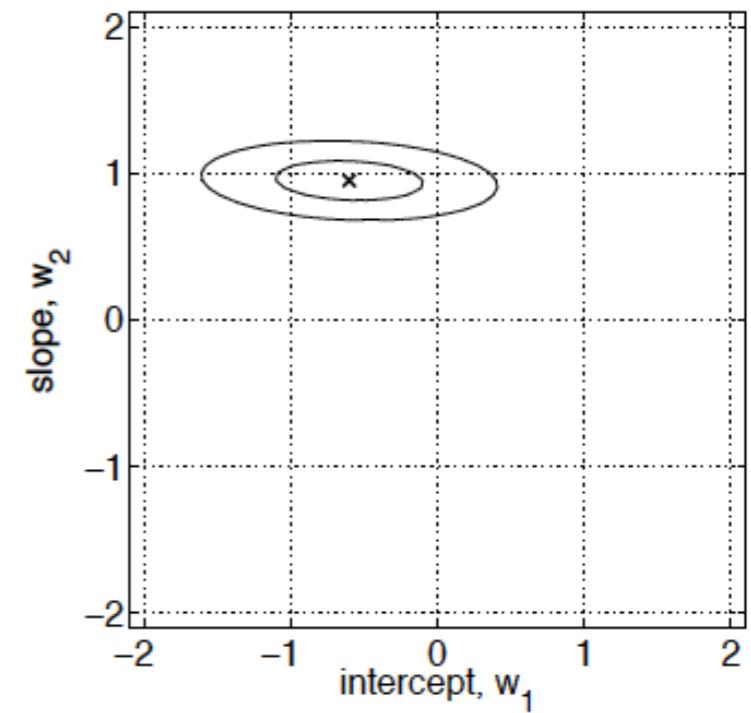
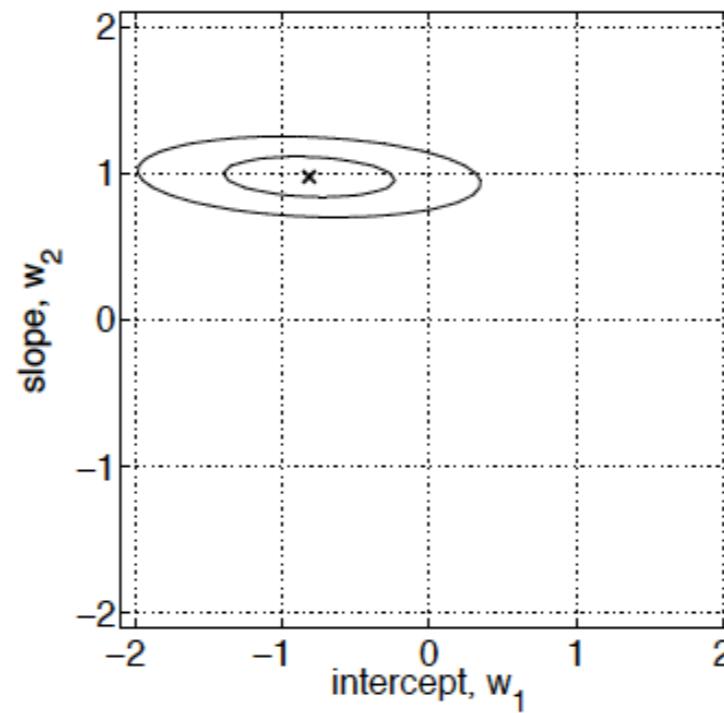
$$p(\mathbf{w}) \sim \mathcal{N}(\bar{\mathbf{0}}, \bar{I}),$$



(a)

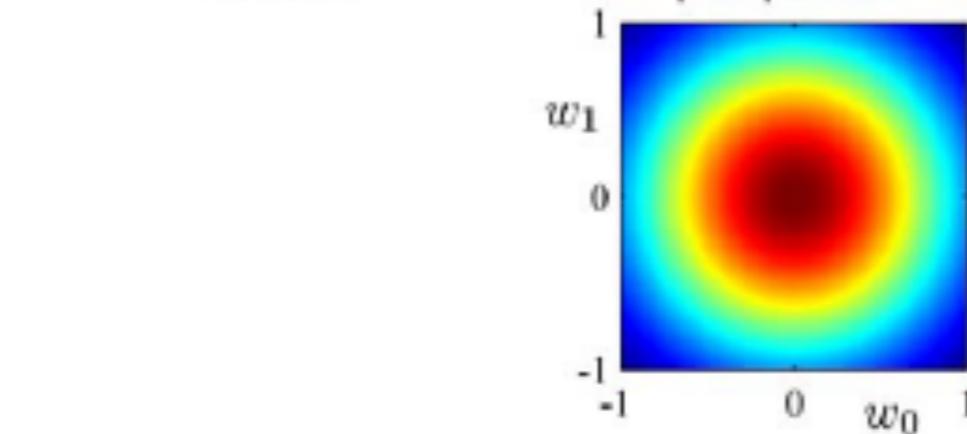


(b)

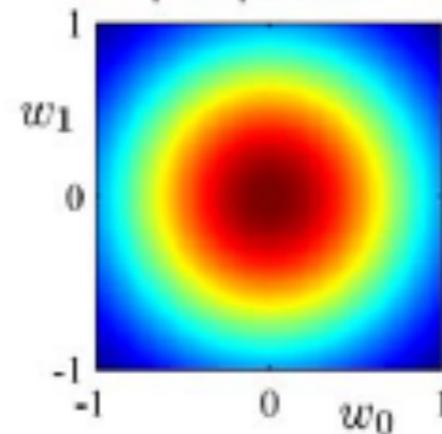


- Sequential update with Bayesian linear regression

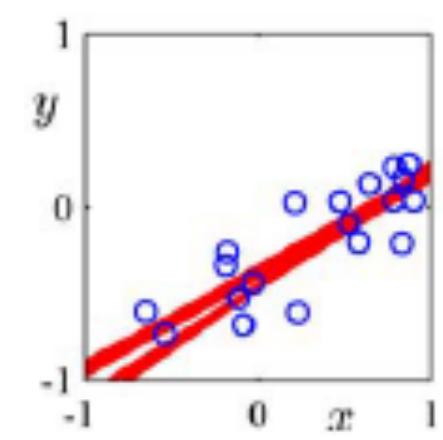
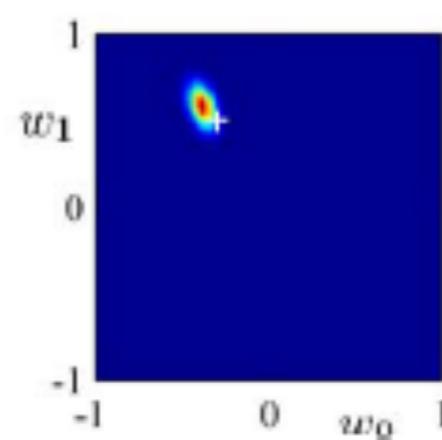
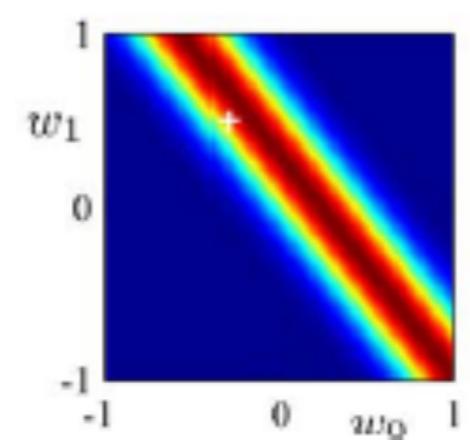
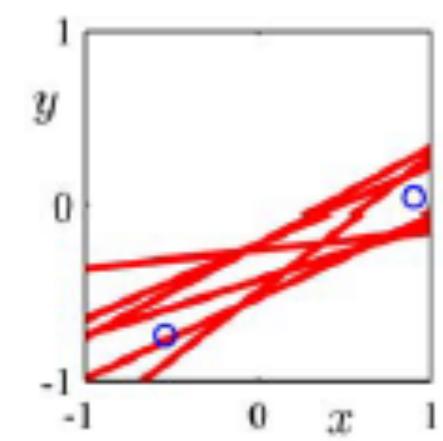
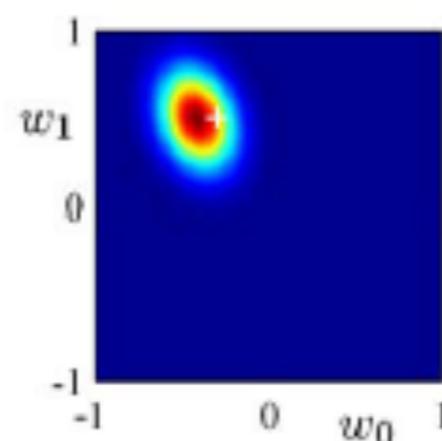
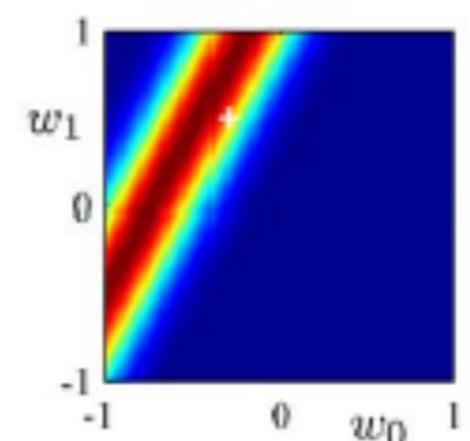
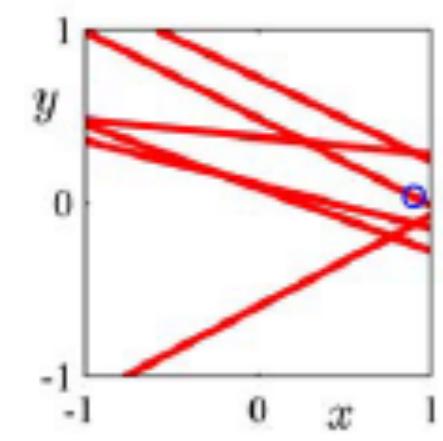
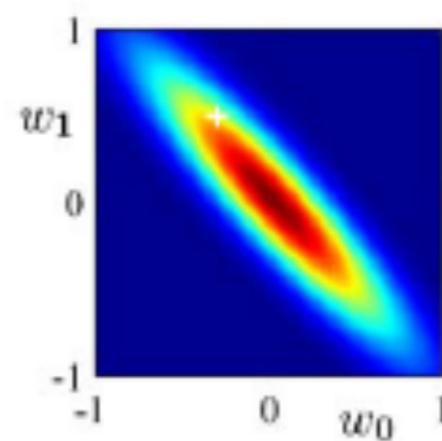
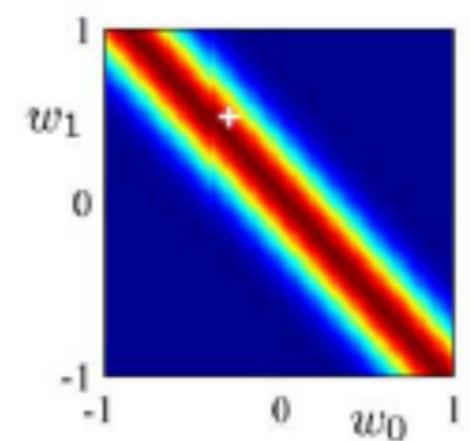
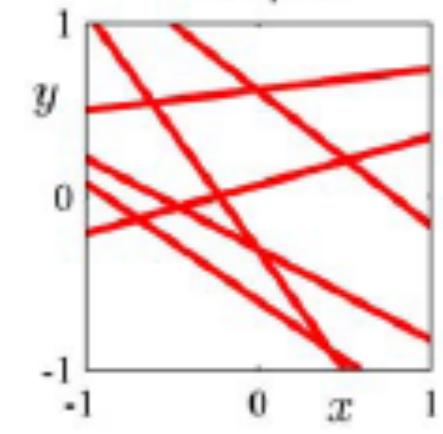
likelihood



prior/posterior



data space



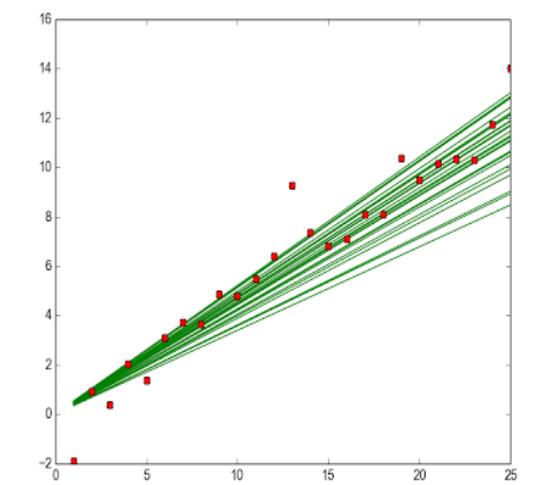
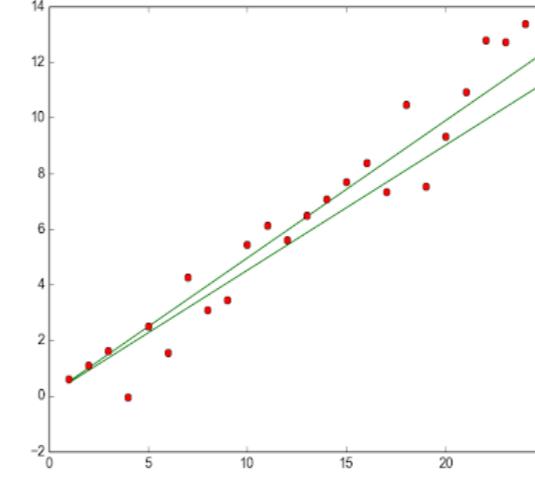
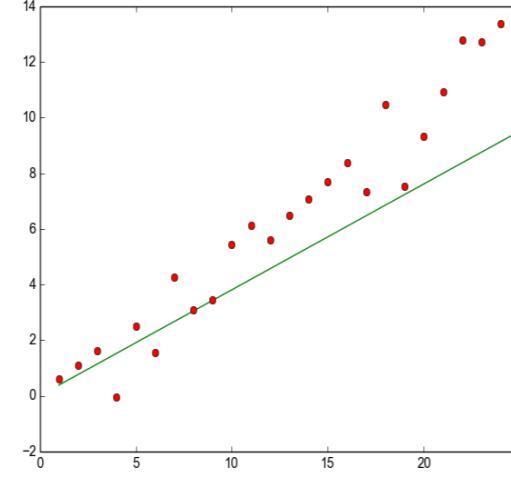
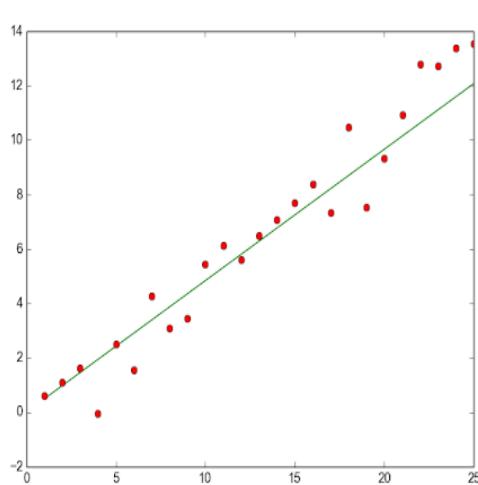
Bayesian linear regression : Prediction

$$\begin{aligned}
 p(f_* | \mathbf{x}_*, X, \mathbf{y}) &= \int p(f_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | X, \mathbf{y}) d\mathbf{w} \\
 &= \mathcal{N}\left(\frac{1}{\sigma_n^2} \mathbf{x}_*^\top A^{-1} X \mathbf{y}, \mathbf{x}_*^\top A^{-1} \mathbf{x}_*\right). \\
 A &= \sigma_n^{-2} X X^\top + \Sigma_p^{-1}.
 \end{aligned}$$

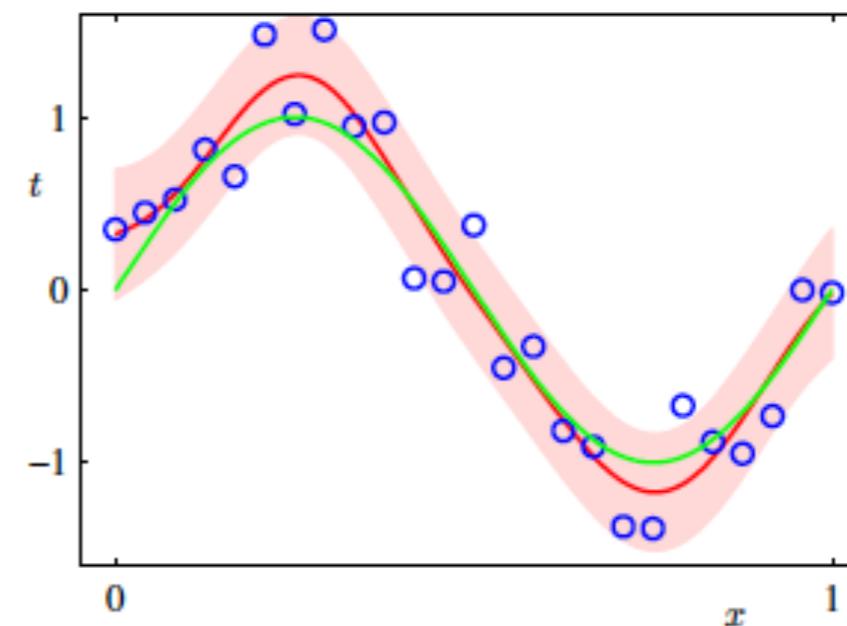
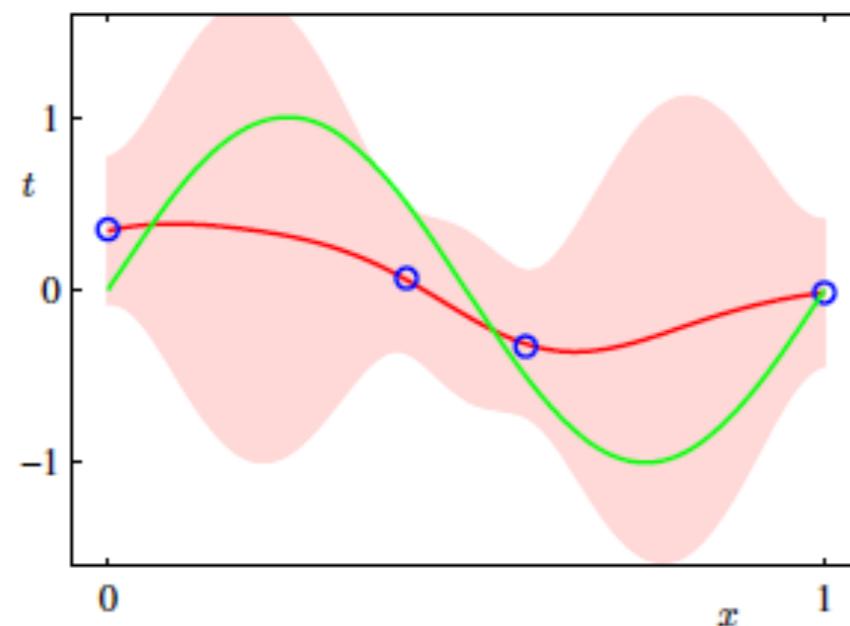
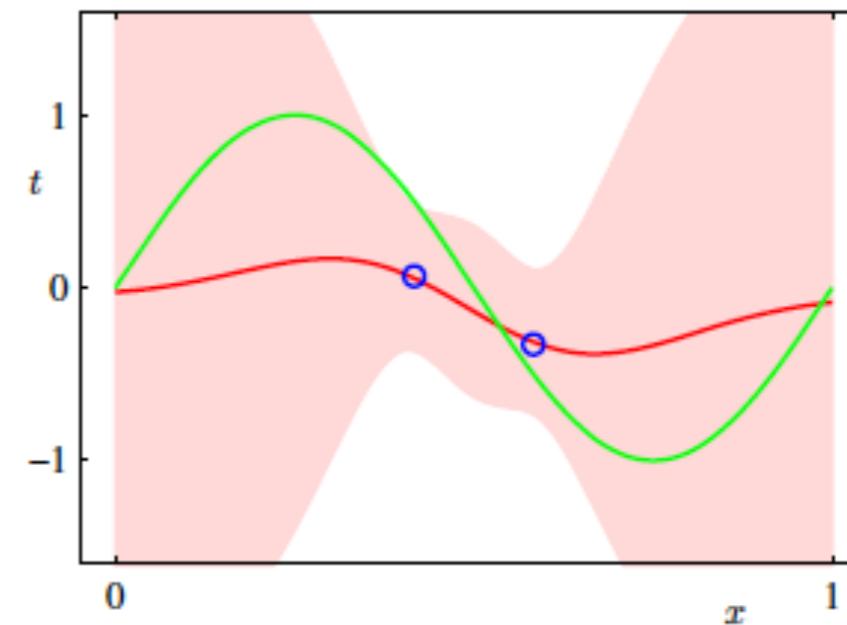
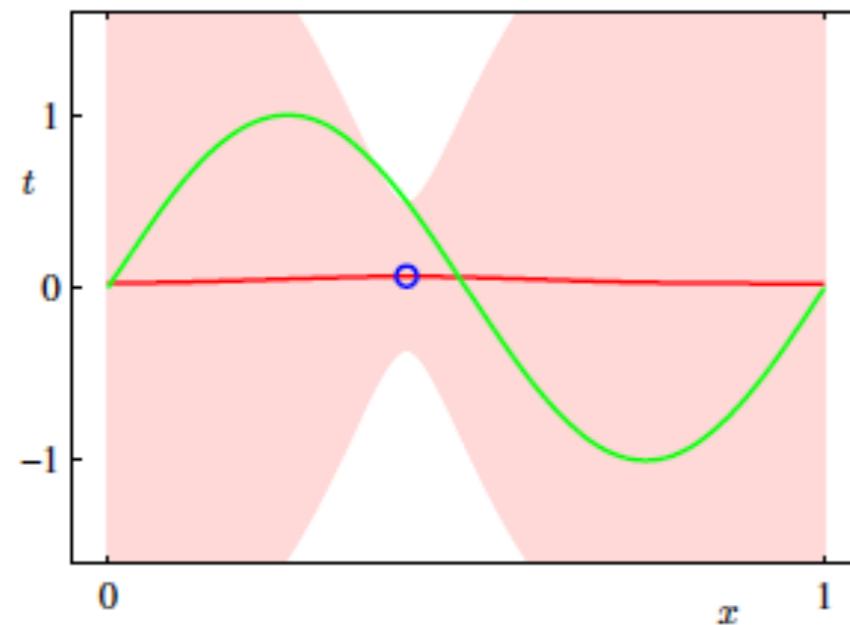
Prediction : $p(y_* | x_*) = \int p(y_* | x_*, w) p(w | D) dw$

$$N(x_*^\top \Sigma_p X (X^\top \Sigma_p X + \sigma_n^2 I)^{-1} y,$$

$$x_*^\top \Sigma_p x_* - x_*^\top \Sigma_p X (X^\top \Sigma_p X + \sigma_n^2 I)^{-1} X^\top \Sigma_p x_*)$$



Bayesian linear regression



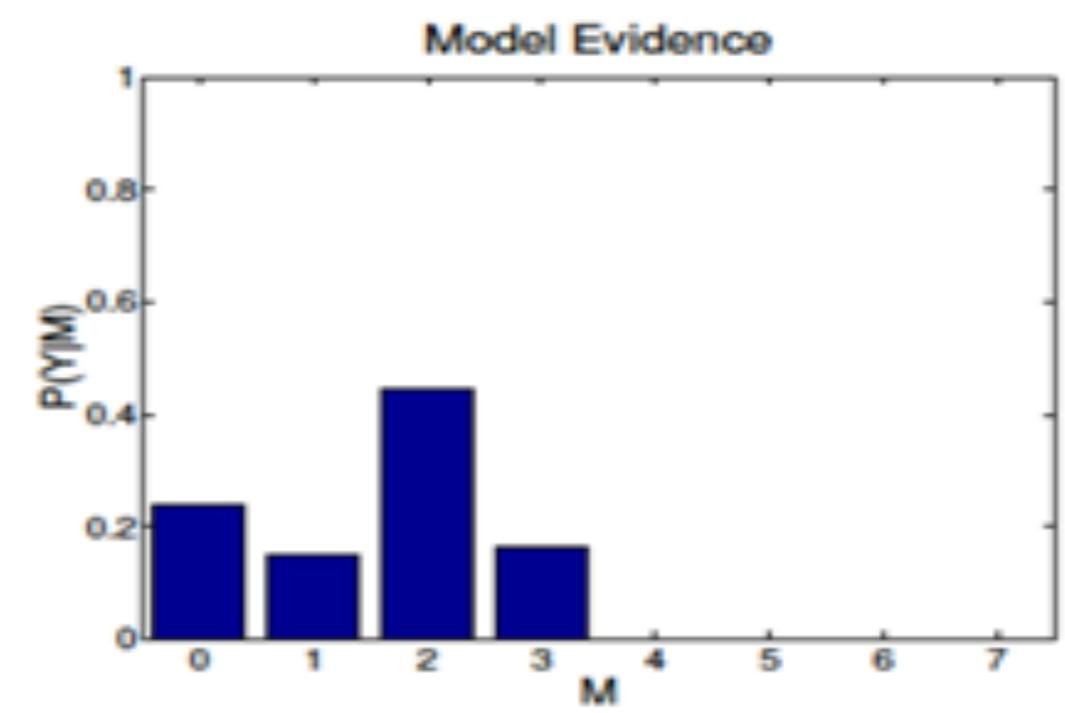
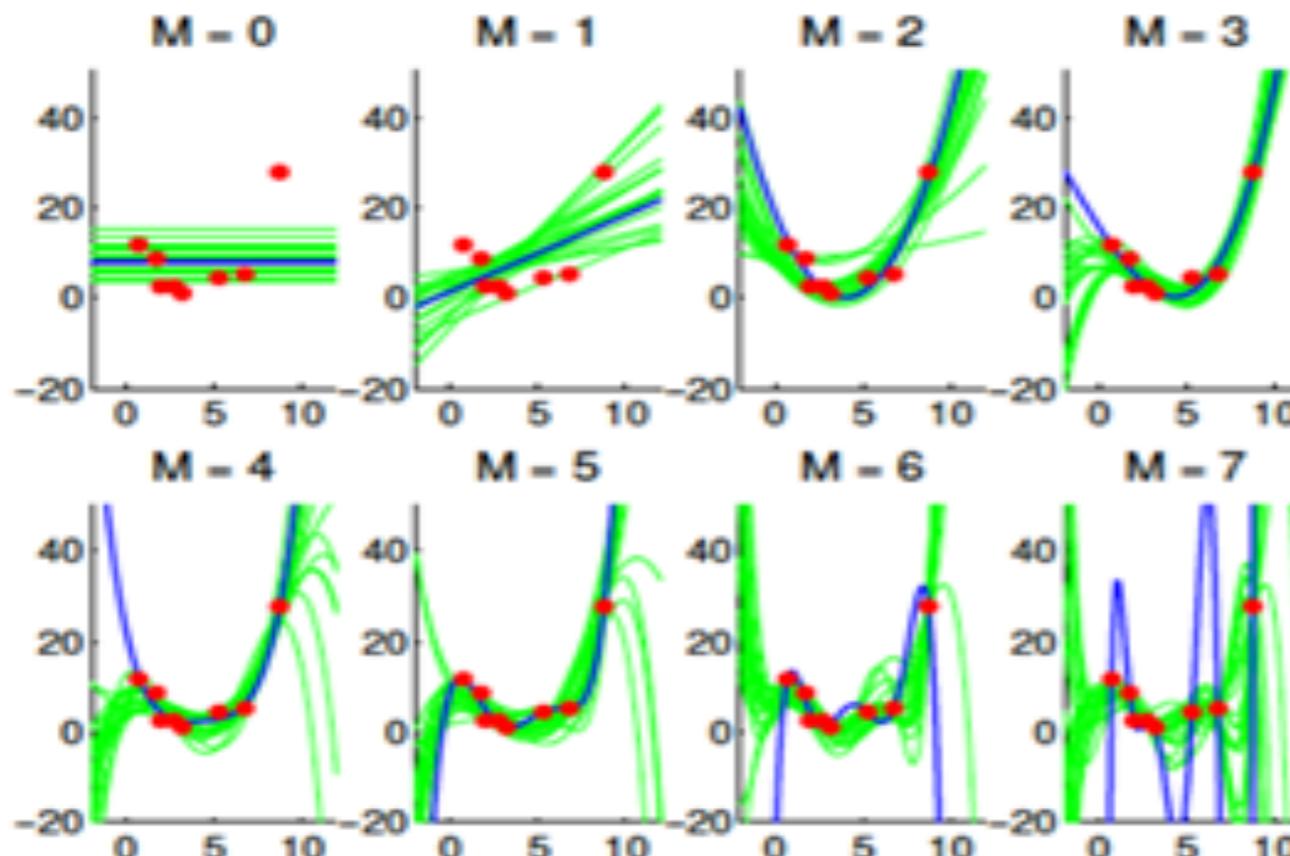
Bayesian linear regression : Model selection

$$p(w|D, M) = \frac{p(D|w, M)p(w|M)}{p(D|M)}$$

ModelSelection : maximize evidence $p(D|M) = \int p(D|w, M)p(w|M)dw$

$$p(y|x, M) = N(y; 0, XX^\top + \sigma^2 I)$$

$$\log p(y|x, M) = -\frac{1}{2}y(XX^\top + \sigma^2 I)^{-1}y - \frac{1}{2}\log |(XX^\top + \sigma^2 I)|$$



Bayesian ML

Treat parameters as random variables to capture model uncertainty
Provides a framework to sequentially update belief about the model/parameters

Bayesian Machine Learning

- Incorporate domain knowledge
 - Through prior
- Modelling uncertainty for safe AI
 - through predictive variance
- Automate machine learning
 - through Evidence
 - Model selection from full data
- Generalize from small data
 - Bayesian model averaging prevents overfitting
- Sequential/online learning
 - Posterior as Prior

Everything follows from two simple rules:

Sum rule: $P(x) = \sum_y P(x, y)$

Product rule: $P(x, y) = P(x)P(y|x)$

$$P(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})}$$

$P(\mathcal{D}|\theta)$ likelihood of θ
 $P(\theta)$ prior probability of θ
 $P(\theta|\mathcal{D})$ posterior of θ given \mathcal{D}

Prediction:

$$P(x|\mathcal{D}, m) = \int P(x|\theta, \mathcal{D}, m)P(\theta|\mathcal{D}, m)d\theta$$

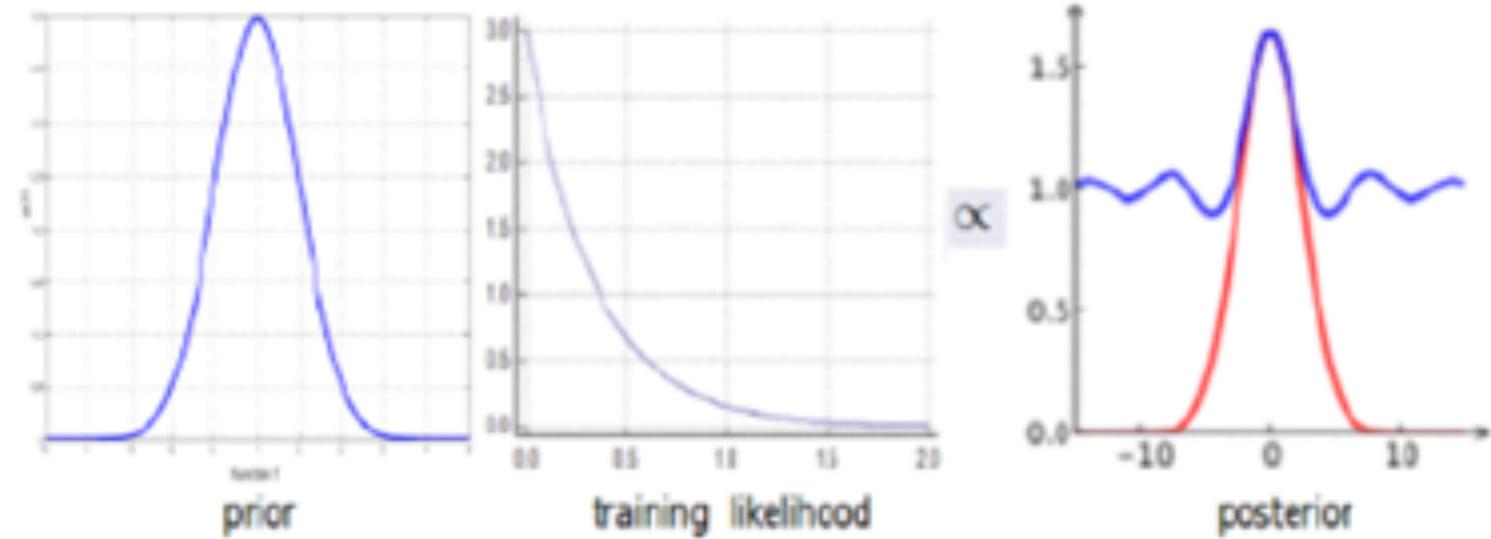
Model Comparison:

$$P(m|\mathcal{D}) = \frac{P(\mathcal{D}|m)P(m)}{P(\mathcal{D})}$$

$$P(\mathcal{D}|m) = \int P(\mathcal{D}|\theta, m)P(\theta|m) d\theta$$

Bayesian machine learning : Challenges

- Choice of prior
- Obtaining the posterior distribution
- What if likelihood and prior are non-conjugates ?
- Bayesian logistic regression
 - Monte-carlo techniques
 - Laplace Approximation
 - Variational Inference
 - Expectation Propagation



References

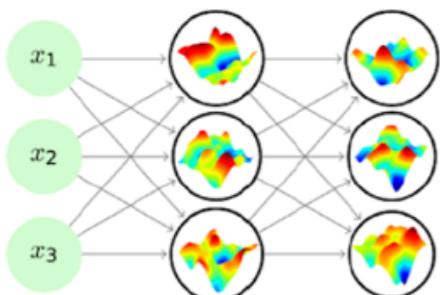
- Pattern Recognition and Machine learning by Christopher Bishop
- Probabilistic machine learning by Kevin Murphy

Bayesian Reasoning And INference (BRAIN)

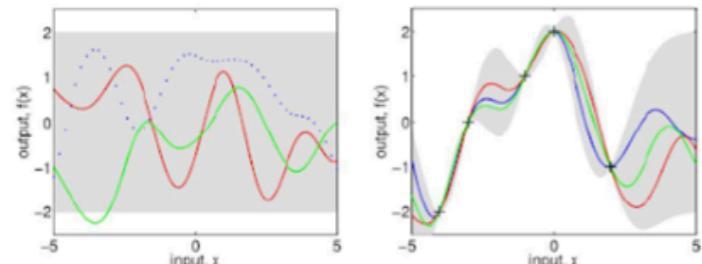
- To understand human learning process and develop machine learning models for artificial intelligence and data science

Models

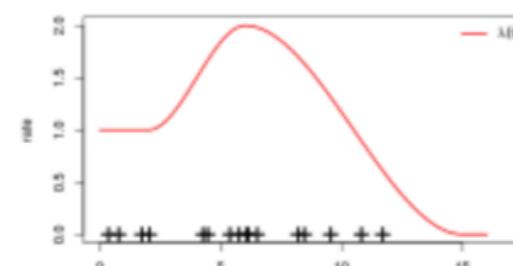
Bayesian deep learning



Bayesian non-parametric



Point Processes



Implicit Generative Model



Applications

Vision and Language



Traffic data analysis



social media analysis



Astrophysical data analysis



<https://sites.google.com/view/brainiith/home>