

M6. Regression: Linear Models -

[a popular (supervised ML) appn. of linear algebra]

Manikandan Narayanan

Week 10 (Sep 29-, 2025)

PRML Jul-Nov 2025 (Grads Section)

Acknowledgment of Sources

- Slides based on content from related
 - Courses:
 - IITM – Profs. Arun/Harish/Chandra’s PRML offerings (slides, quizzes, notes, etc.), Prof. Ravi’s “Intro to ML” slides – cited respectively as [AR], [HR]/[HG], [CC], [BR] in the bottom right of a slide.
 - India – NPTEL PR course by IISc Prof. PS. Sastry (slides, etc.) – cited as [PSS] in the bottom right of a slide.
 - Books:
 - PRML by **Bishop**. (content, figures, slides, etc.) – cited as **[CMB]**
 - Pattern Classification by Duda, Hart and Stork. (content, figures, etc.) – **[DHS]**
 - Mathematics for ML by Deisenroth, Faisal and Ong. (content, figures, etc.) – **[DFO]**
 - Information Theory, Inference and Learning Algorithms by David JC MacKay – **[DJM]**

Outline for Module M5

- M6. Regression (Linear models)
 - M6.0 Introduction/context
 - Problem formulation
 - Motivating examples, incl. basis functions
 - M6.1 Linear regression approaches
 - Direct approach I: least-squares (error) (w_{LS})
 - Direct approach II: least-squares (error) with regularization (w_{RLS})
 - Discriminative model approach I: MLE ($w_{ML} = w_{LS}$)
 - Discriminative model approach IIa: Bayesian linear regression ($w_{MAP} = w_{RLS}$)
 - Discriminative model approach IIb: Bayesian linear regression (posterior & predictions)
 - M6.2 Model Complexity/Selection
 - Motivation (hyperparameter tuning to avoid overfitting)
 - Frequentist view (bias-variance decomposition)
 - M6.3 Concluding thoughts

Outline for Module M5

- M6. Regression (Linear models)
 - M6.0 Introduction/context
 - Problem formulation
 - Motivating examples, incl. basis functions
 - M6.1 Linear regression approaches
 - Direct approach I: least-squares (error) (w_{LS})
 - Direct approach II: least-squares (error) with regularization (w_{RLS})
 - Discriminative model approach I: MLE ($w_{ML} = w_{LS}$)
 - Discriminative model approach IIa: Bayesian linear regression ($w_{MAP} = w_{RLS}$)
 - Discriminative model approach IIb: Bayesian linear regression (posterior & predictions)
 - M6.2 Model Complexity/Selection
 - Motivation (hyperparameter tuning to avoid overfitting)
 - Frequentist view (bias-variance decomposition)
 - M6.3 Concluding thoughts

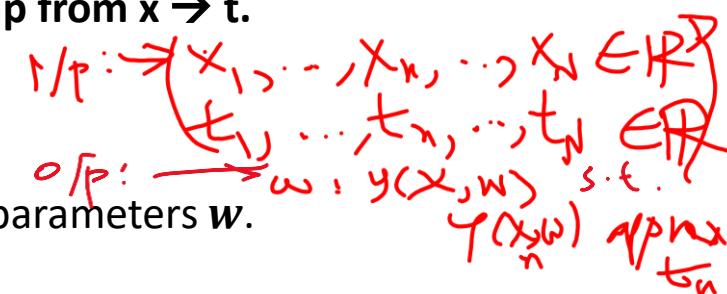
Context: ML Paradigms

- Unsupervised Learning (informally aka “learning patterns from (unlabelled) data”)
 - Density estimation
 - Clustering
 - Dimensionality reduction
- Supervised Learning (informally aka curve-fitting or function approximation or “function learning from (labelled) data”)
 - Learn an input and output map (features x to target t)
 - **Regression:** continuous output/target t
 - Classification: categorical output

Regression: the problem

- Learn a map from input variables (features x) to a continuous output variable (target t).
Informally, known as function approximation/learning or curve fitting, since
Given $\{x_n, t_n\}_{n=1\dots N}$ pairs, we seek a function $y_w(x)$ that “approximates” the map from $x \rightarrow t$.

Linear regression assumes $y_w(x) := y(x, w)$ is a **linear function of** the adjustable parameters w .
It could be linear or non-linear in x .



- A foundational supervised learning problem/algorithm:
 - Practical limitations for complex data, but sets analytical foundation for other sophisticated learning algorithms.
 - Due to its simplicity, first and predominant choice of statistical model in many applied areas with moderate sample sizes:
(E.g., in bioinformatics: to adjust for known confounding factors (covariates) in Disease Genomics and Genome-wide Association (GWAS) Studies, Causal inference such as in Mendelian Randomization, etc.)
- Our approach in this lecture: Mostly [CBM, Chapters 1,3].

Regression: what does “approximate” mean? recall three approaches (from Decision Theory)

- Generative model approach:
 - (I) Model $p(t, \mathbf{x})$
 - (I) Infer $p(t|\mathbf{x})$ from $p(t, \mathbf{x})$
 - (D) Take conditional mean/median/mode/any other optimal decision outcome as $y(\mathbf{x})$

- Discriminative model approach:
 - (I) Model $p(t|\mathbf{x})$ directly
 - (D) Take conditional mean/median/mode/any other optimal decision outcome as $y(\mathbf{x})$

- Direct regression approach:
 - (D) Learn a regression function $y(\mathbf{x})$ directly from training data

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

Linear regression: from linear combination of input variables ($\mathbf{x} \in \mathbb{R}^D$) to that of basis functions ($\phi(\mathbf{x}) \in \mathbb{R}^M$)

- Simplest model of linear regn. involving D input vars.:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_{(1)} + \dots + w_D x_{(D)}$$

$$= w_0 \cdot 1 + \sum_{j=1}^D w_j x_j = [\mathbf{w}_0 \ \mathbf{w}_1 \ \dots \ \mathbf{w}_D] \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} = \mathbf{w}^T \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$$

- Model of linear regn. involving M basis fns. (**fixed** non-linear fns. of the input vars.):

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 \phi_1(\mathbf{x}) + \dots + w_{M-1} \phi_{M-1}(\mathbf{x})$$

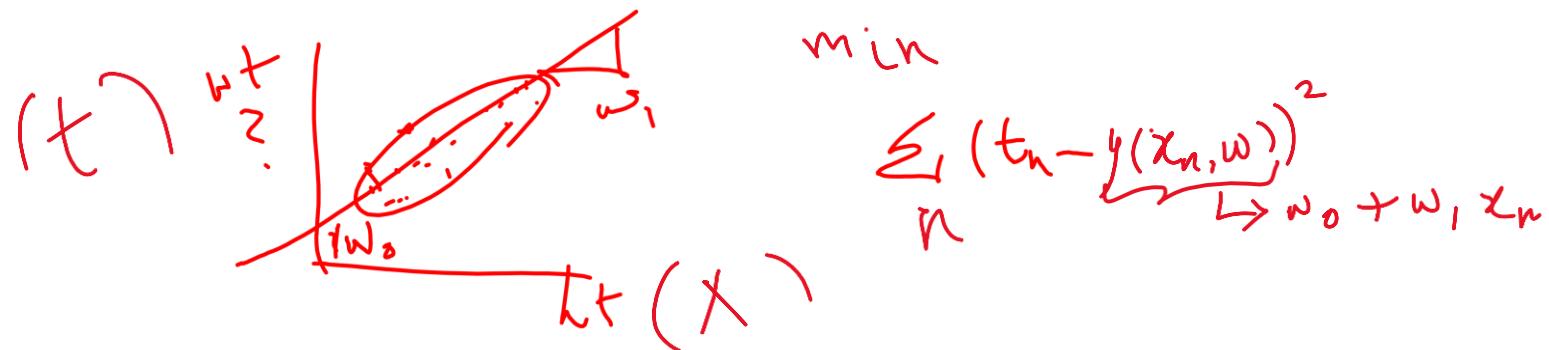
$$= w_0 \cdot 1 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

$$= \mathbf{w}^T \phi(\mathbf{x})$$

($\phi: \mathbb{R}^D \rightarrow \mathbb{R}^M$, with convention $\phi_0(\mathbf{x}) = 1$)

Linear regression: recall standard examples

- Predicting weight t from height x : $y(x, w)$ linear in both x and w .



- Estimation of fetal weight t (actually $\log_{10} t$) from ultrasound measurements: $y(x, w)$ linear in w , not x .

Author	Components	Formula
Ferrero	AC, FL	$10^{(0.77125 + 0.13244 \cdot AC - 0.12996 \cdot FL - 1.73588 \cdot AC^2/1,000 + 3.09212 \cdot FL \cdot AC/1,000 + 2.18984 \cdot FL/AC)}$ (g, cm)
Hadlock I	BPD, HC, AC, FL	$10^{(1.3596 + 0.0064 \cdot HC + 0.0424 \cdot AC + 0.174 \cdot FL + 0.00061 \cdot BPD \cdot AC - 0.00386 \cdot AC \cdot FL)}$ (g, cm)
Hadlock III	BPD, AC, FL	$10^{(1.335 - 0.0034 \cdot AC \cdot FL + 0.0316 \cdot BPD + 0.0457 \cdot AC + 0.1623 \cdot FL)}$ (g, cm)
Hadlock V	BPD, AC	$10^{(1.1134 + 0.05845 \cdot AC - 0.000604 \cdot AC^2 - 0.007365 \cdot BPD^2 + 0.000595 \cdot BPD \cdot AC + 0.1694 \cdot BPD)}$ (g, cm)
Shepard	BPD, AC	$10^{(-1.7492 + 0.166 \cdot BPD + 0.046 \cdot AC - 0.002546 \cdot AC \cdot BPD)}$ (kg, cm)

[From Hoopmann et al. *Fetal Diagn Ther.* 2011;30(1):29-34.
doi:10.1159/000323586]

More examples of lin. regn. of basis functions

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_{(1)} + \dots + w_D x_{(D)}$$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_{(1)} + w_2 \sqrt{x_{(2)}} + w_3 x_{(1)} \sqrt{x_{(2)}} + w_4 \log(x_{(1)})$$

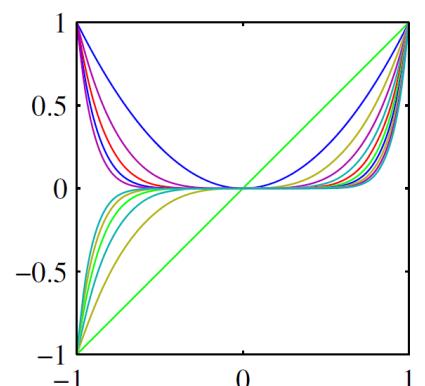
$$\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^5$$

$$\mathbf{x} = (x_{(1)}, x_{(2)}) \rightarrow (1, x_{(1)}, \sqrt{x_{(2)}}, x_{(1)}\sqrt{x_{(2)}}, \dots)$$

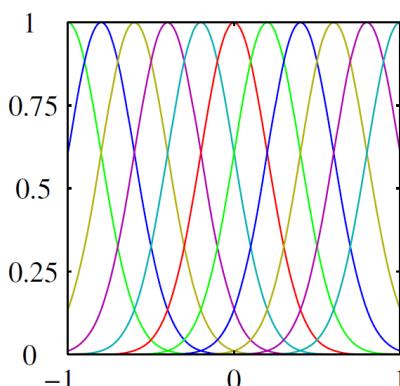
$$\phi: \mathbb{R}^D \rightarrow (\mathbb{R}^{D+1})$$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

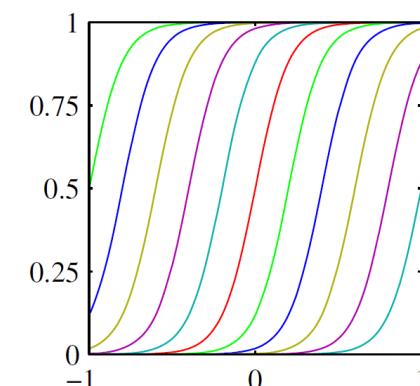
$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$



Polynomials (& extn.
to spatially restricted
splines)



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



$$\phi_j(x) = \sigma \left(\frac{x - \mu_j}{s} \right)$$

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

$(X \in \mathbb{R}^D)$
 $\underbrace{\quad}_{n} \quad \underbrace{\quad}_{M}$
 $\phi(x_n) \in \mathbb{R}$

Fourier basis
(sinusoidal fns.),
Wavelets, etc.

[CMB]

Outline for Module M5

- M6. Regression (Linear models)
 - M6.0 Introduction/context
 - Problem formulation
 - Motivating examples, incl. basis functions
 - **M6.1 Linear regression approaches**
 - Direct approach I: least-squares (error) (w_{LS})
 - Direct approach II: least-squares (error) with regularization (w_{RLS})
 - Discriminative model approach I: MLE ($w_{ML} = w_{LS}$)
 - Discriminative model approach IIa: Bayesian linear regression ($w_{MAP} = w_{RLS}$)
 - Discriminative model approach IIb: Bayesian linear regression (posterior & predictions)
 - M6.2 Model Complexity/Selection
 - Motivation (hyperparameter tuning to avoid overfitting)
 - Frequentist view (bias-variance decomposition)
 - M6.3 Concluding thoughts

Linear regression: a direct approach

Approach: minimize sum-of-squares error; aka
least-squares solution/approach

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

where $y(x, \mathbf{w}) = \mathbf{w}^T \phi(x)$

$$\mathcal{D}_N = \left\{ \begin{bmatrix} x_n \\ t_n \end{bmatrix} \right\}_{n=1}^N$$

↓
 $\phi(x_n) \in \mathbb{R}^M$

Solution: w_{LS} that minimizes $E_D(w)$ (via matrix notation)

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2. = \frac{1}{2} \|\vec{\phi}w - \vec{t}\|^2$$

→ basis
↑ data pts

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \cdots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \cdots & \phi_{M-1}(x_N) \end{pmatrix} \quad N \times M$$

Normal equations from setting gradient to zero, using $N \times M$ design matrix:

$$\phi^T \phi w_{LS} = \phi^T \vec{t} \quad \textcircled{1}$$

$$\Rightarrow w_{LS} = (\phi^T \phi)^{-1} \cdot \phi^T \vec{t}$$

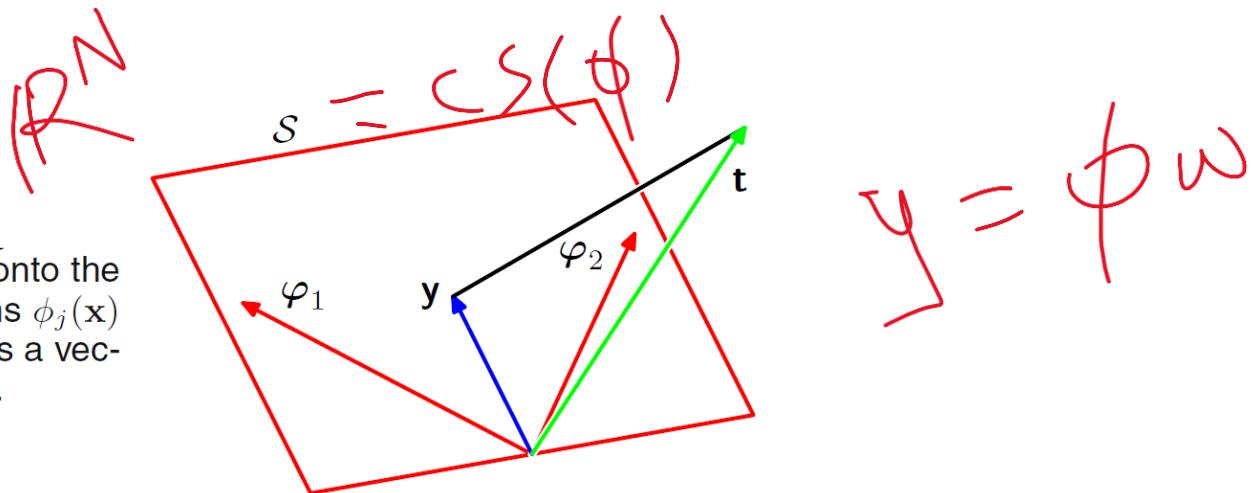
- (a) Prove (1) always has a soln.
- (b) If ϕ has linearly indep. cols, then $\phi^T \phi$ is invertible.
- (c) What if $\phi^T \phi$ is not invertible?

$$\vec{t} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{pmatrix}$$

$$\vec{y} = \phi w$$

Soln.: Geometry of (least-squares) sol.

α orthogonal projection of the data vector t onto the subspace spanned by the basis functions $\phi_j(x)$ in which each basis function is viewed as a vector φ_j of length N with elements $\phi_j(x_n)$.



[CMB, Figure 3.2]

$$\text{min}_{w \in \mathbb{R}^m} \|(\phi w - t)\|^2 \quad \text{--- (A)}$$

$$\text{min}_{y \in CS(\phi)} \|y - t\|^2 \quad \text{--- (B)}$$

[CMB, HR]

LA Refresher

- See Appendix for LA refresher on
 - related topic of $Ax=b$ when no solution is possible, and
 - some matrix-vector gradient formulas/tricks.

Recall LA: To solve $Ax = b$, we premult. by A^T , and simply solve $A^T A x = A^T b$.

Linear Equation Solving

$$Ax = b$$

$\textcircled{1} \rightarrow$ No solution if $b \notin \text{column space}(A)$

\rightarrow Unique solution if $b \in \text{column space}(A)$ & A has lin. ind. columns.

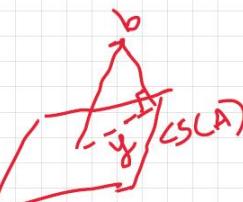
\rightarrow Infinite solutions if $b \in \text{CS}(A)$ & A has lin. dep. columns.

$$\begin{bmatrix} 1 & 3 \\ 2 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$\textcircled{2}$ Least-squares soln possible (related to lin. regn)

$$\min_x \|Ax - b\|^2 = \min_{y \in \text{CS}(A)} \|y - b\|^2$$

$(Ax^* = y^*)$ (y^* is proj. of b onto $\text{CS}(A)$)



$$\min_{x \in \mathbb{R}^n} [(Ax - b)^T (Ax - b)] \rightarrow f(x)$$

$$\nabla f(x) = 2A^T(Ax - b) \stackrel{\text{set to } 0}{=} 0$$

$$\Rightarrow A^T(Ax^* - b) = 0$$

$$\Rightarrow A^T A x^* = A^T b$$

$$\Rightarrow x^* = (A^T A)^{-1} A^T b$$

(normal eqn.)

Ex.: Prove:

- 1) at least one soln. x^* exists for the normal eqn.
- 2) soln. x^* unique if $(A^T A)$ is invertible ($\Leftrightarrow A$ has lin. indep. cols.)
- 3) infinite solns. x^* if $(A^T A)$ is non-invertible ($\Leftrightarrow A$ has lin. dep. cols.)

Ex.:

i) Prove $NS(A) = NS(A^T A)$.

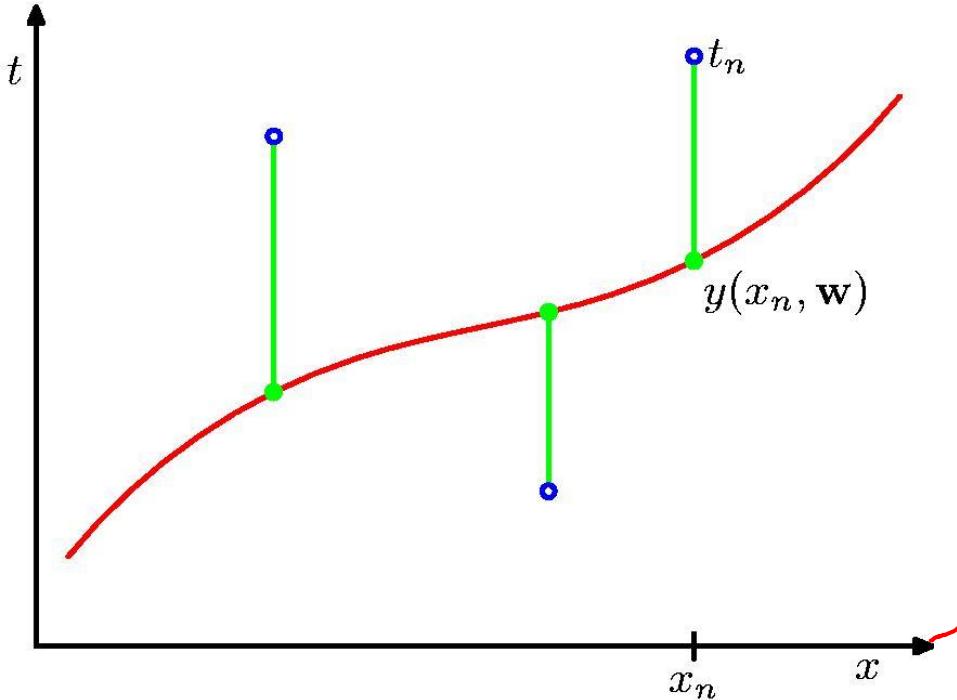
ii) Use orthog. complementarity of $NS(A^T)$, $CS(A)$ to derive normal eqns.

Outline for Module M5

- M6. Regression (Linear models)
 - M6.0 Introduction/context
 - Problem formulation
 - Motivating examples, incl. basis functions
 - **M6.1 Linear regression approaches**
 - Direct approach I: least-squares (error) (w_{LS})
 - **Direct approach II: least-squares (error) with regularization (w_{RLS})**
 - Discriminative model approach I: MLE ($w_{ML} = w_{LS}$)
 - Discriminative model approach IIa: Bayesian linear regression ($w_{MAP} = w_{RLS}$)
 - Discriminative model approach IIb: Bayesian linear regression (posterior & predictions)
 - M6.2 Model Complexity/Selection
 - Motivation (hyperparameter tuning to avoid overfitting)
 - Frequentist view (bias-variance decomposition)
 - M6.3 Concluding thoughts

From sum-of-squares to regularized error!

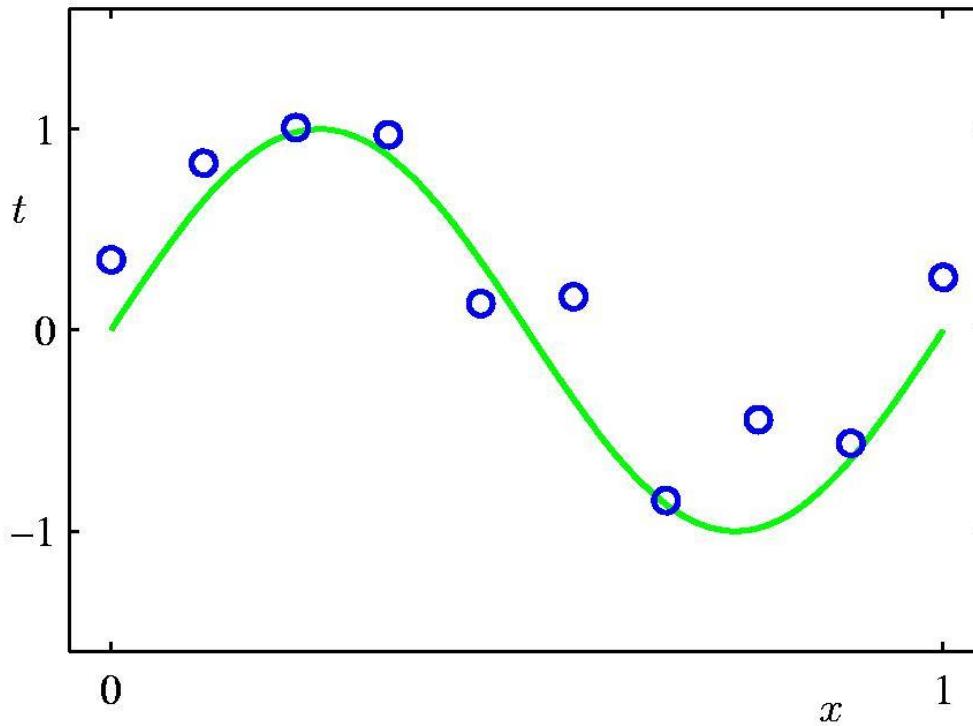
Running example: polynomial curve fitting ((via sum-of-squares error function / least squares approach))



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

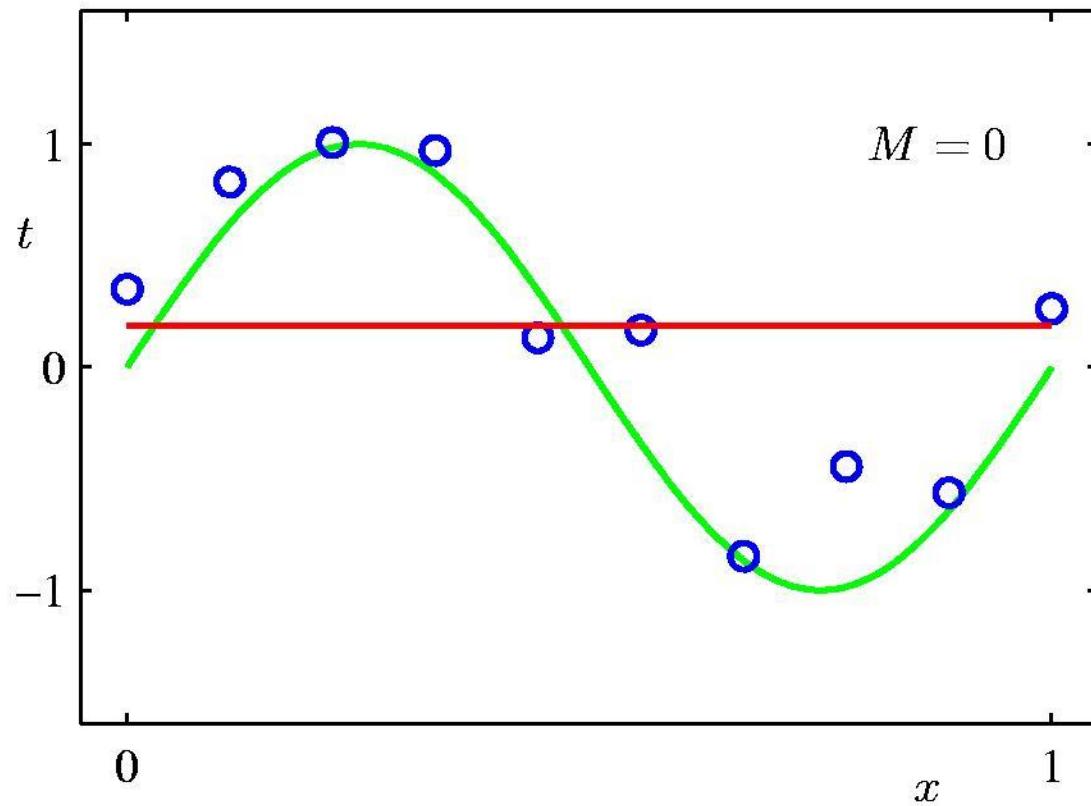
$$\text{where } y(x, \mathbf{w}) = \mathbf{w}^T \phi(x)$$

Polynomial Curve Fitting

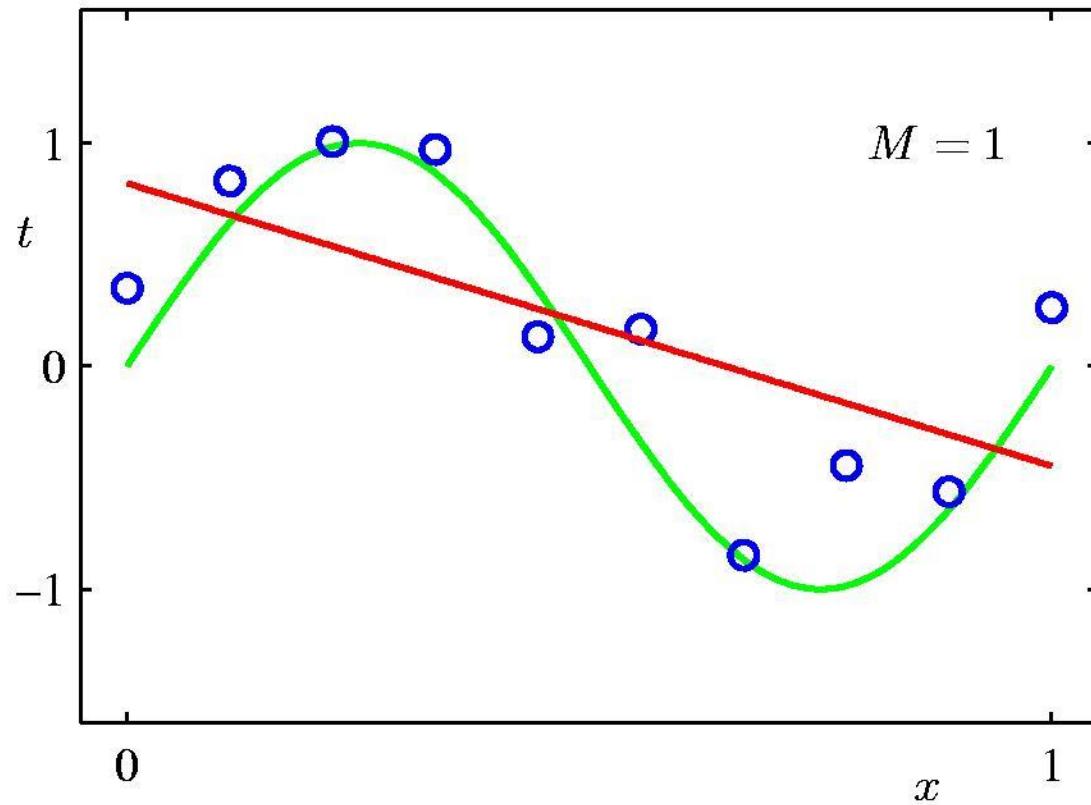


$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

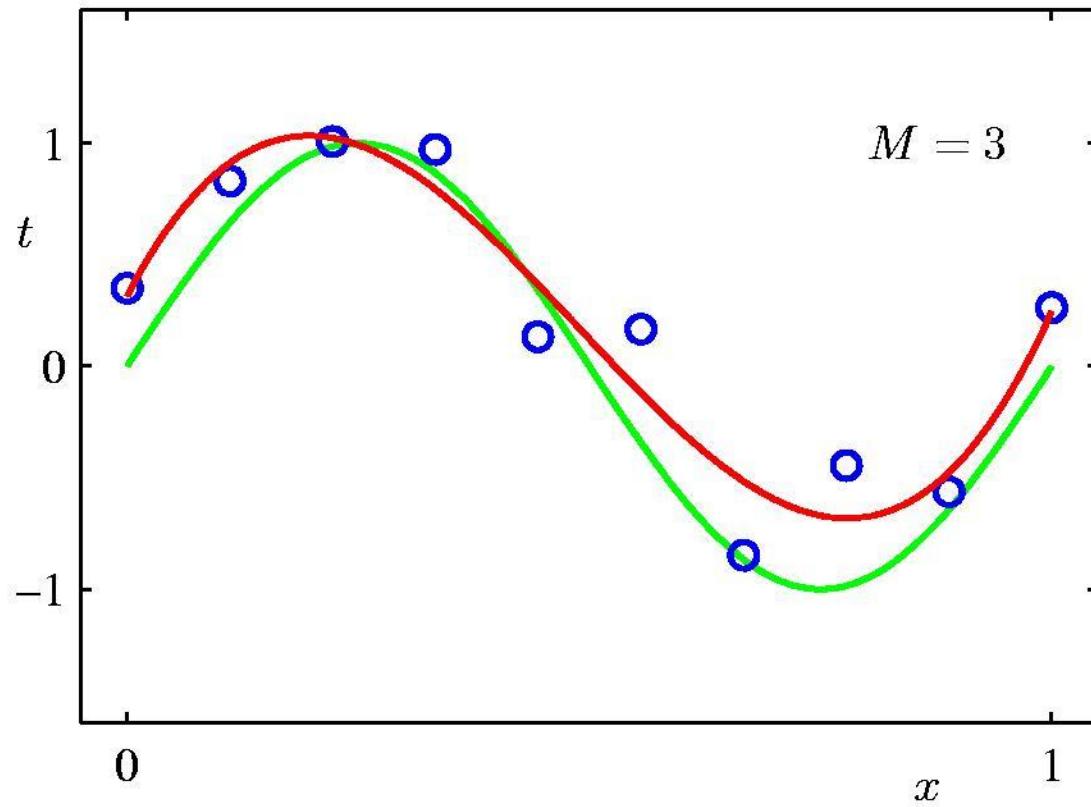
0th Order Polynomial



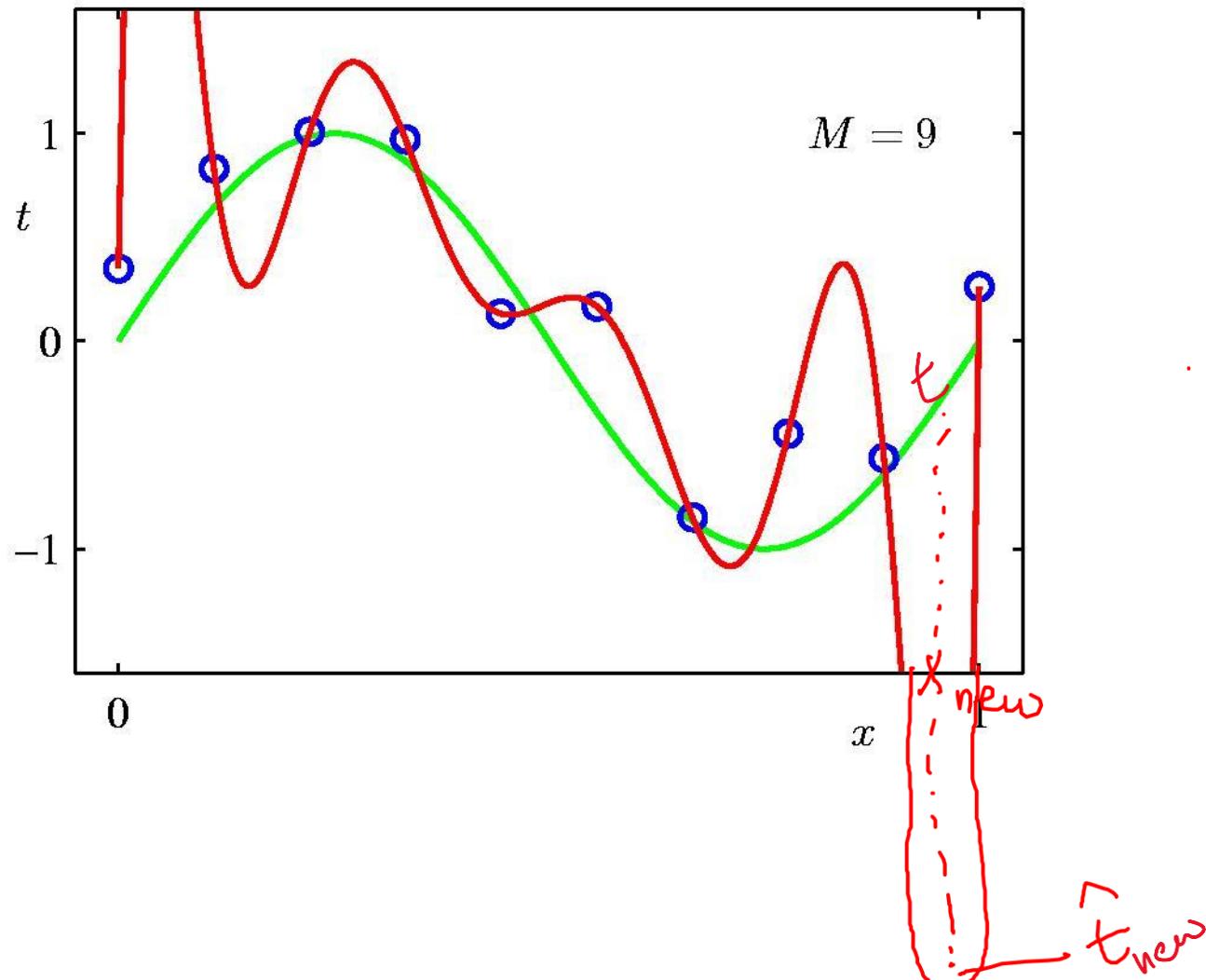
1st Order Polynomial



3rd Order Polynomial



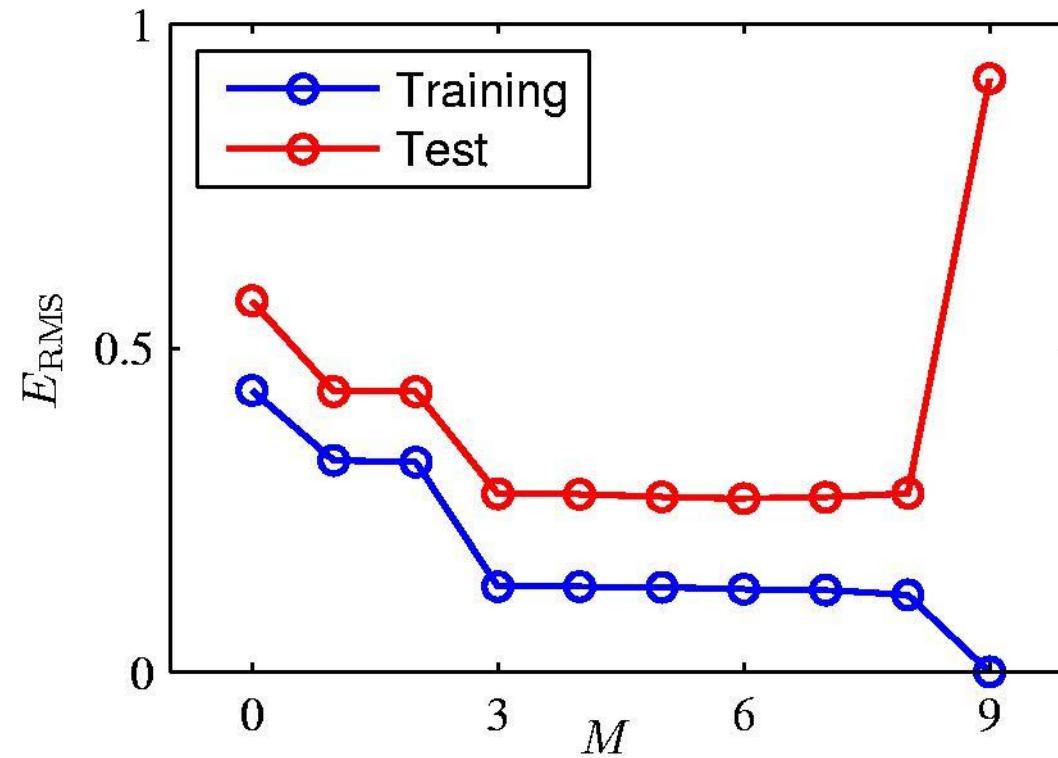
9th Order Polynomial



[CMB]

Brainstorm: What degree polynomial would you choose?

Over-fitting



Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

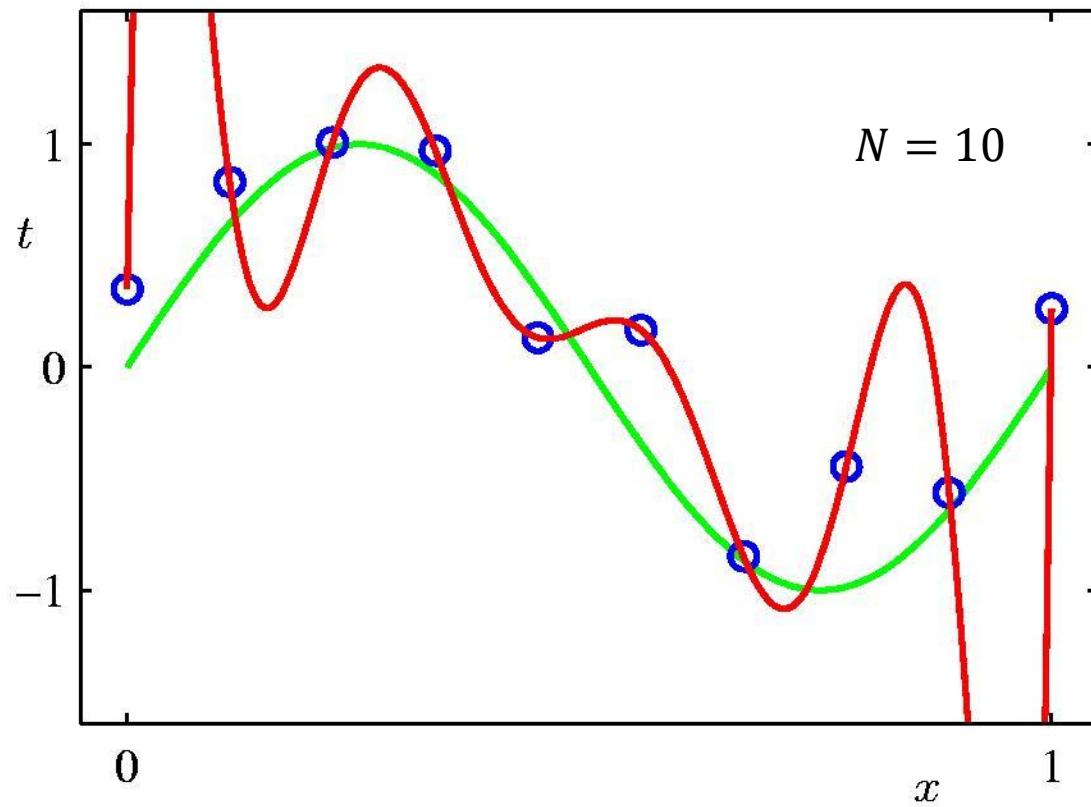
Polynomial Coefficients

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	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

The role of data set size N?

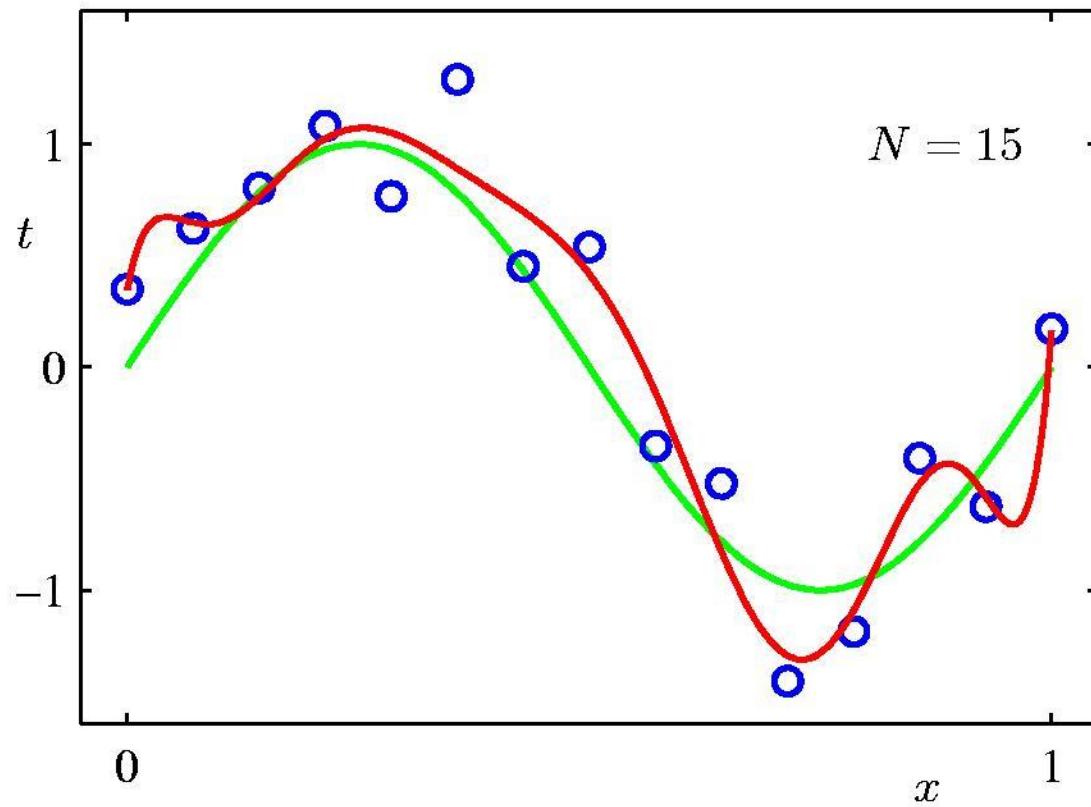
Data Set Size: $N = 10$

9th Order Polynomial



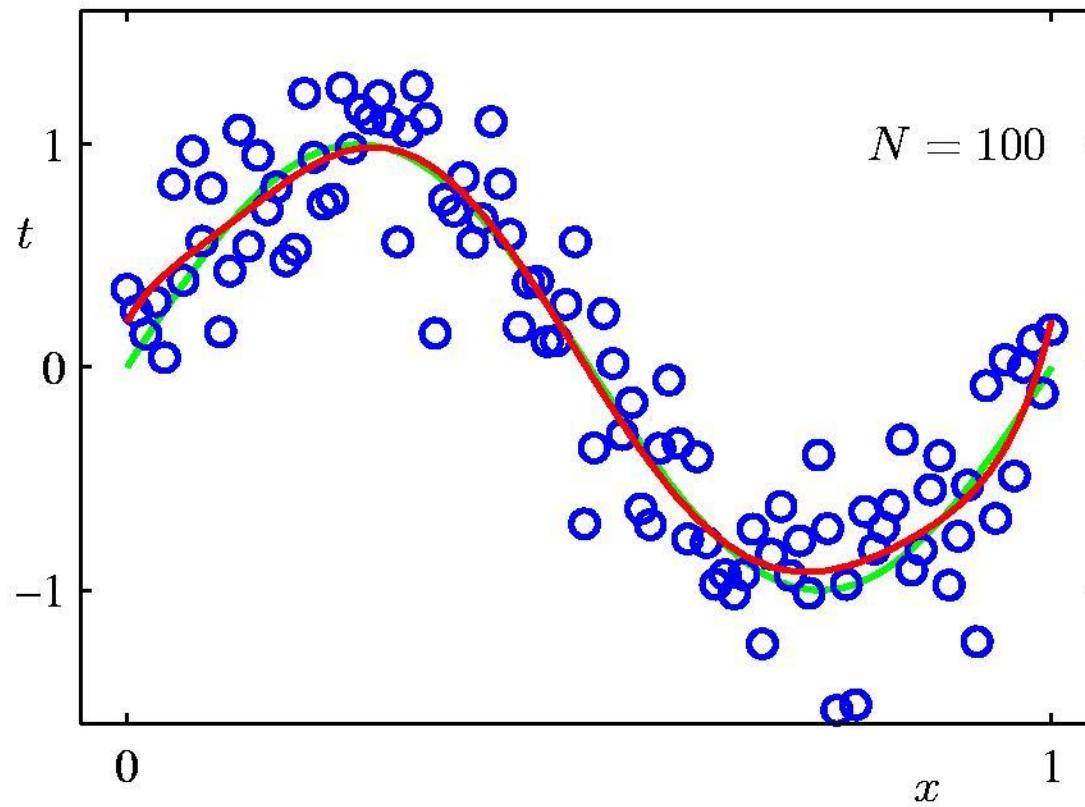
Data Set Size: $N = 15$

9th Order Polynomial



Data Set Size: $N = 100$

9th Order Polynomial



How to take care of both data set size and model complexity tradeoffs?

Regularization

- Penalize large coefficient values

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



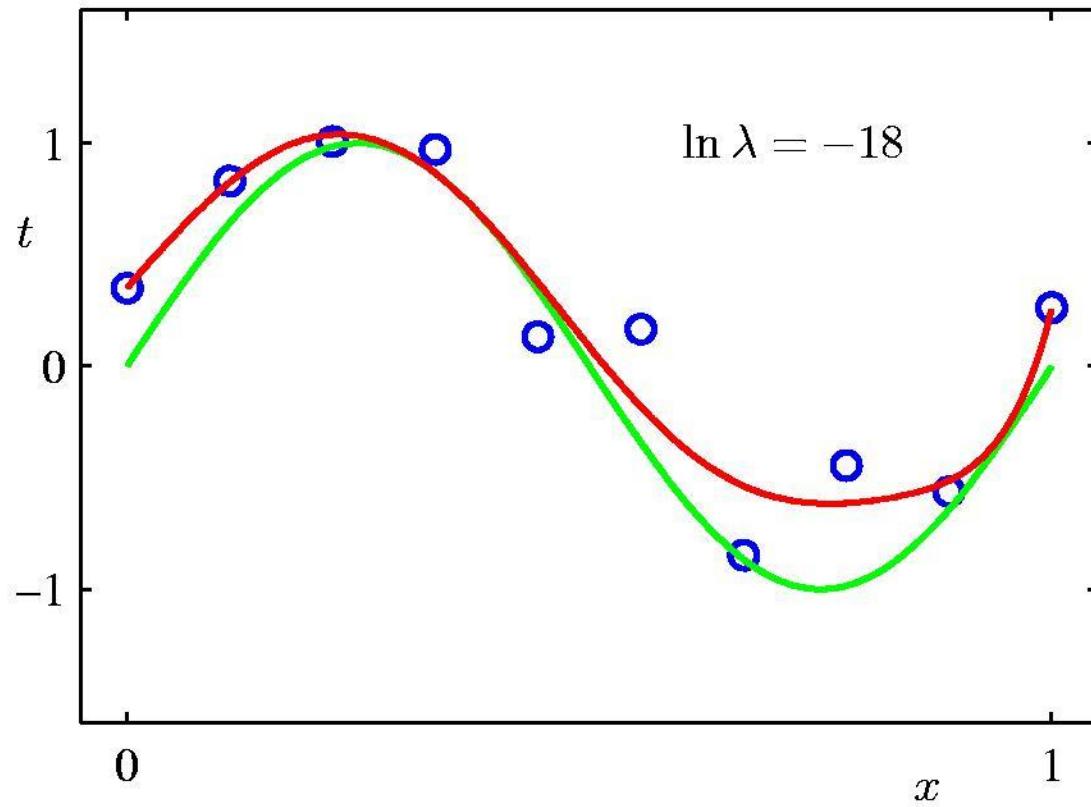
favors complex models penalizes complex models

Polynomial Coefficients

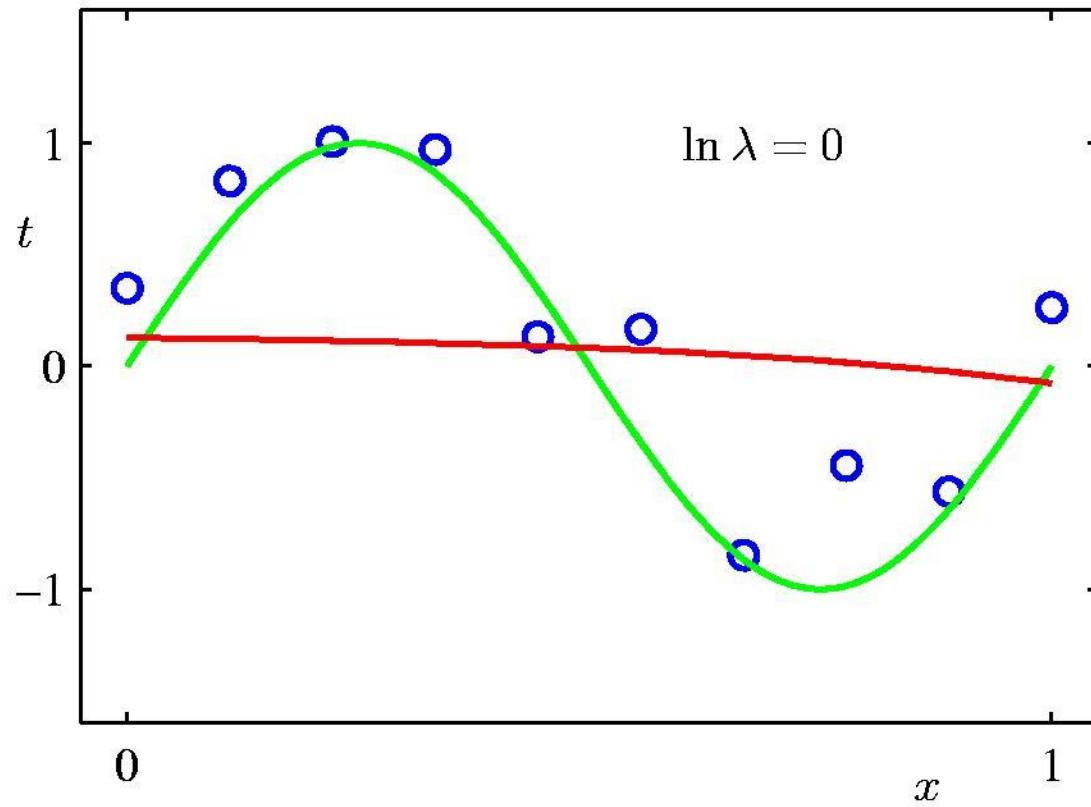
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53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

Regularization:

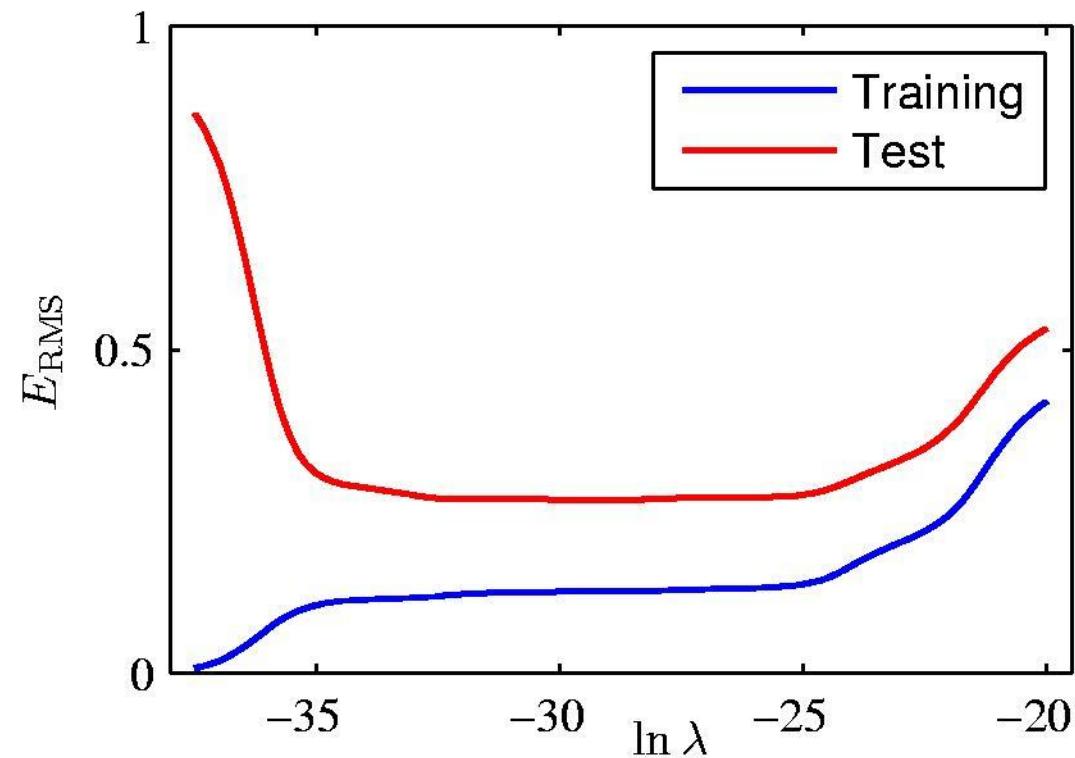
$$\ln \lambda = -18$$



Regularization: $\ln \lambda = 0$



Regularization: E_{RMS} vs. $\ln \lambda$



more complex model ← → less complex model

[CMB]

Now, let's see how to solve the minimization problem!

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

Matrix notation gradient formula (chain rule product “builds towards the left”):

$$\begin{aligned} f(x) &= g(Ax) \\ \nabla f(x) &= A^T \nabla g(Ax) \end{aligned} \quad \begin{aligned} f: \mathbb{R}^d &\rightarrow \mathbb{R}, \quad g: \mathbb{R}^n \rightarrow \mathbb{R} \\ A &\in \mathbb{R}^{n \times d} \end{aligned}$$

Regularized Least Squares (1): Solution for ridge regression

- Consider the error function:

$$\tilde{E}(w) = E_D(w) + \lambda E_W(w)$$

Data term + Regularization term

λ is called the regularization coefficient.

- With the sum-of-squares error function and a quadratic regularizer, we get

$$\frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2 + \frac{\lambda}{2} w^T w$$

$\frac{1}{2} \| \phi w - t \|^2$

$\left(\frac{\lambda}{2} \| w \|^2 \right)$

- which is minimized by

$$w_{RLS} = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T t.$$

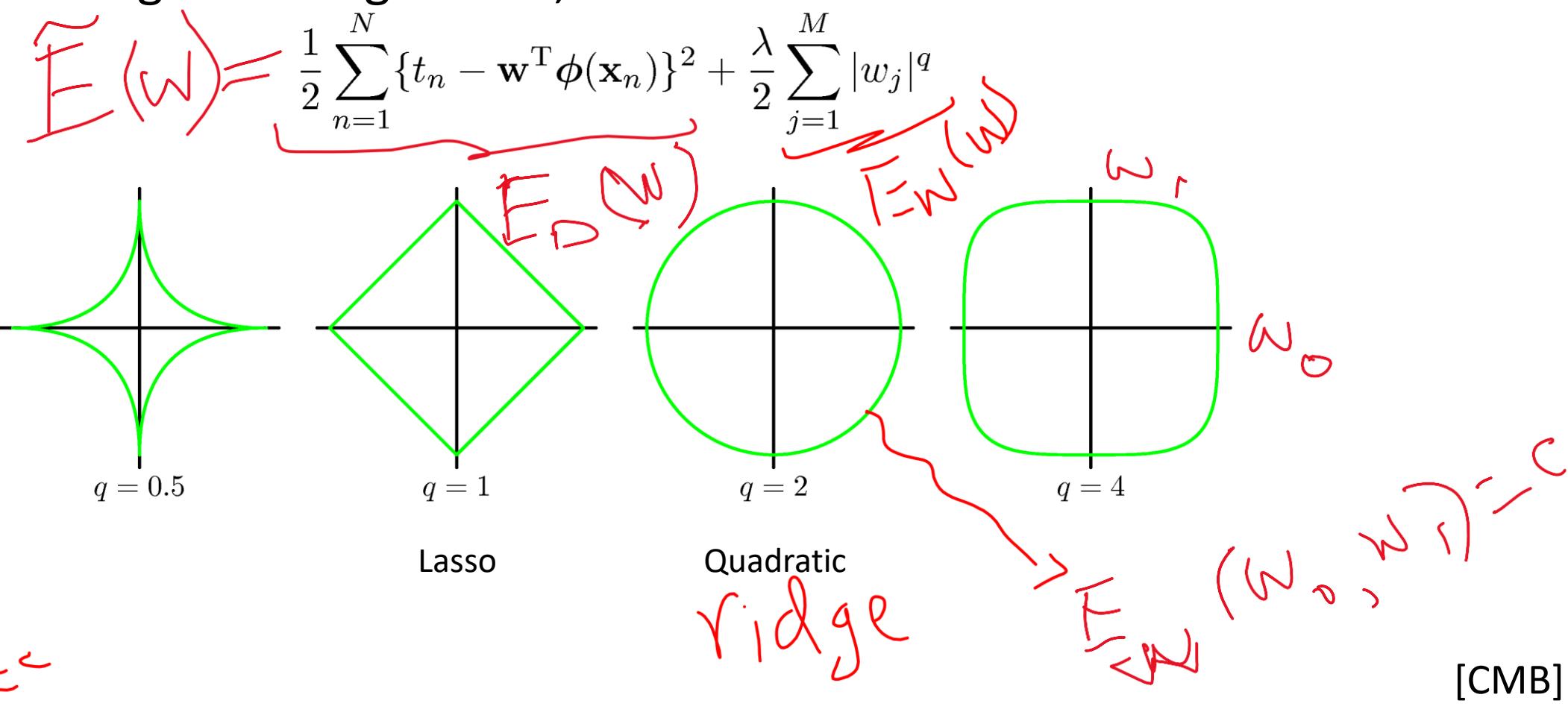
$M \times M$

Ex: ST $\lambda I + \Phi^T \Phi$ is invertible for $\lambda > 0$

$$\phi \in \mathbb{R}^{N \times M}$$

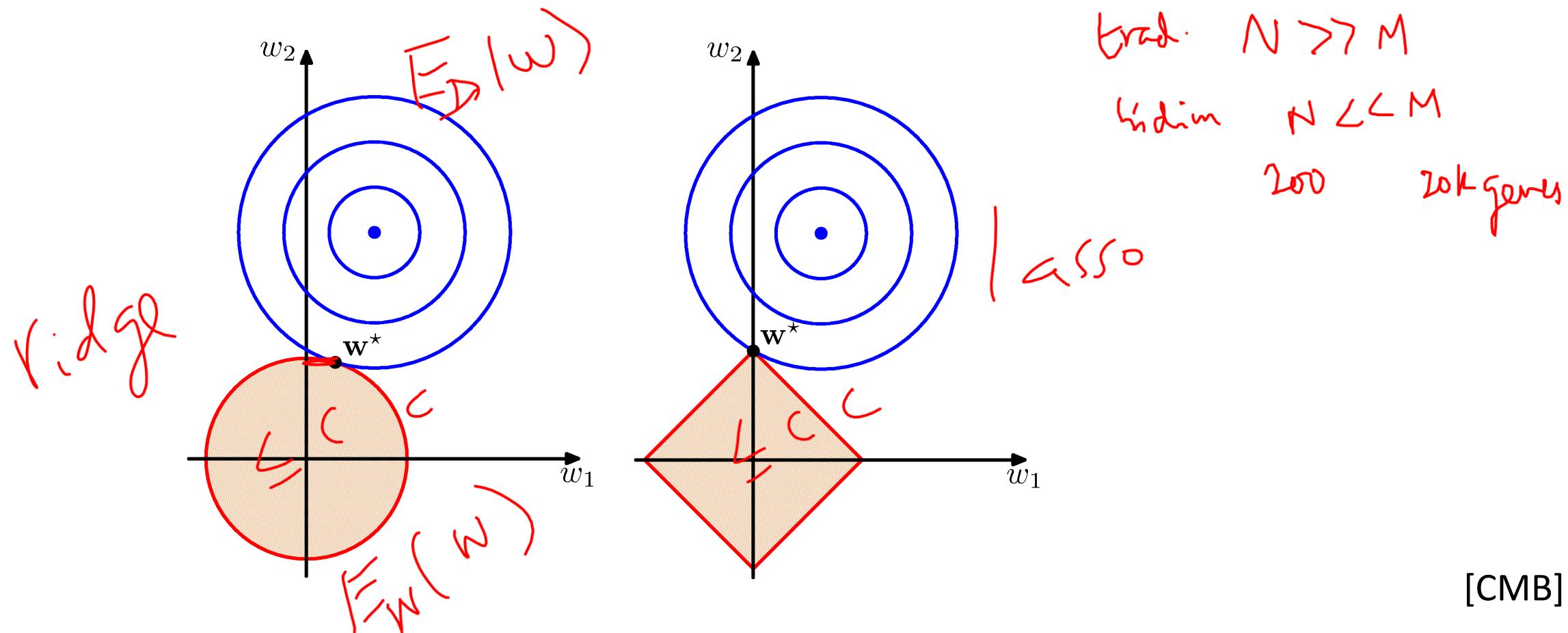
Regularized Least Squares (2)

- With a more general regularizer, we have



Regularized Least Squares (3)

- Lasso tends to generate sparser solutions than a ridge (quadratic) regularizer.
- Regularization aka penalization/weight-decay in ML or parameter shrinkage in statistics literature.



Outline for Module M5

- M6. Regression (Linear models)
 - M6.0 Introduction/context
 - Problem formulation
 - Motivating examples, incl. basis functions
 - **M6.1 Linear regression approaches**
 - Direct approach I: least-squares (error) (w_{LS})
 - Direct approach II: least-squares (error) with regularization (w_{RLS})
 - **Discriminative model approach I: MLE ($w_{ML} = w_{LS}$)**
 - Discriminative model approach IIa: Bayesian linear regression ($w_{MAP} = w_{RLS}$)
 - Discriminative model approach IIb: Bayesian linear regression (posterior & predictions)
 - M6.2 Model Complexity/Selection
 - Motivation (hyperparameter tuning to avoid overfitting)
 - Frequentist view (bias-variance decomposition)
 - M6.3 Concluding thoughts

A different view of min. $E(w)$ or its regularized version?

Q: How do you convert this intuition/empirical-art into science, and derive $E(w)$ or its (other) regularized versions systematically?

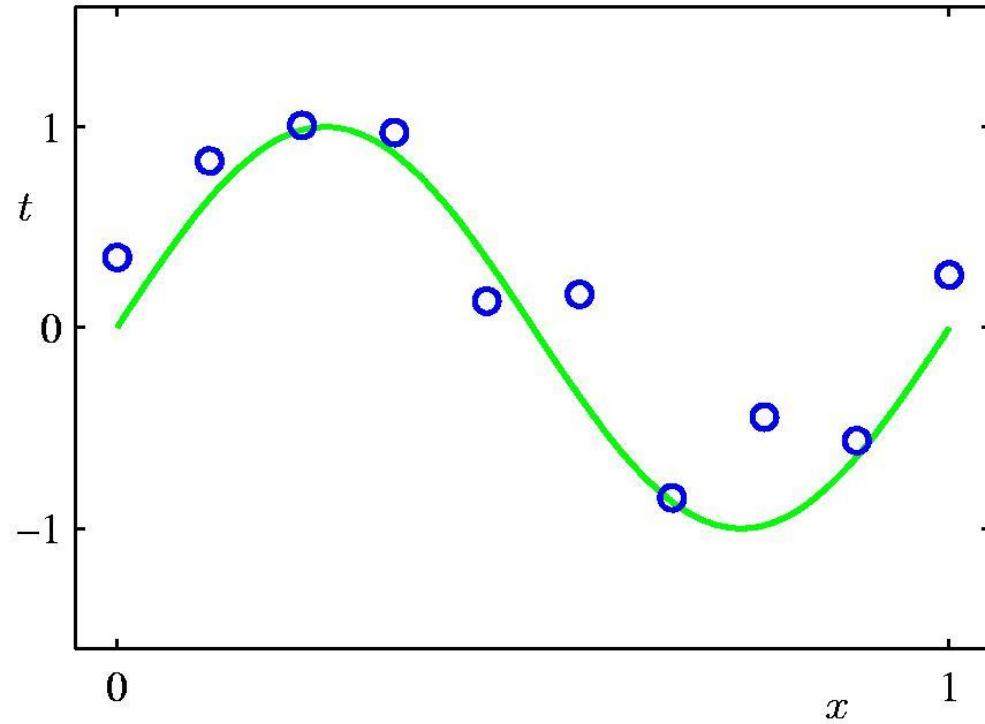
A: Probabilistic view helps. Discriminative Approach: Model $p(t|x)$

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

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

Q: What has regression got to do with our previous topic: density estimation?

- Brainstorm: how to model $P(t|x)$?



Q: What has regression got to do with our previous topic: density estimation?

- Brainstorm: how to model $P(t|x)$?

$$p(t) \quad p_{M_\theta}(t)$$

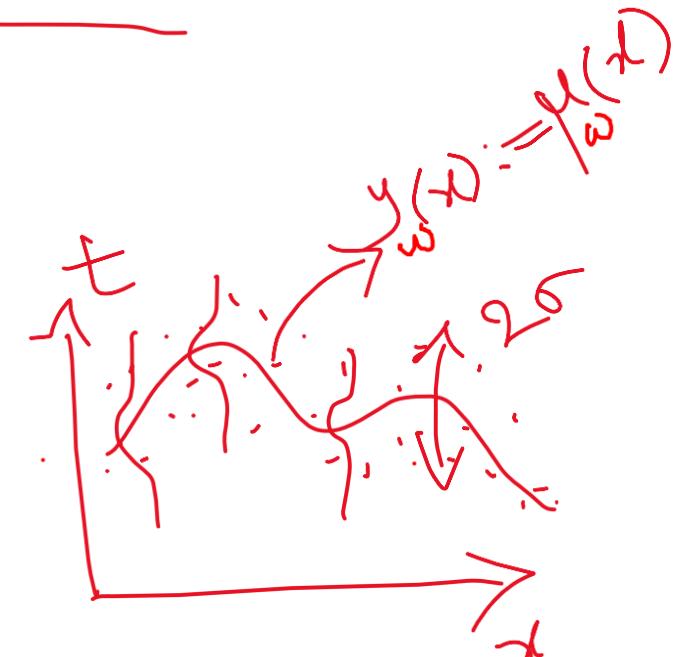
$$p(t|x) \quad p_{M_{\theta(x)}}(t)$$

E.g.)

$$p(t) \approx N_{\mu, \sigma^2}(t) dt$$



$$p(t|x) \approx N_{\mu(x), \sigma^2}(t) dt$$



Q: What has regression got to do with our previous topics?

A: $P(t|x)$ captures the input-output map. Steps involved are:

(1) Model/estimate $P(t|x)$

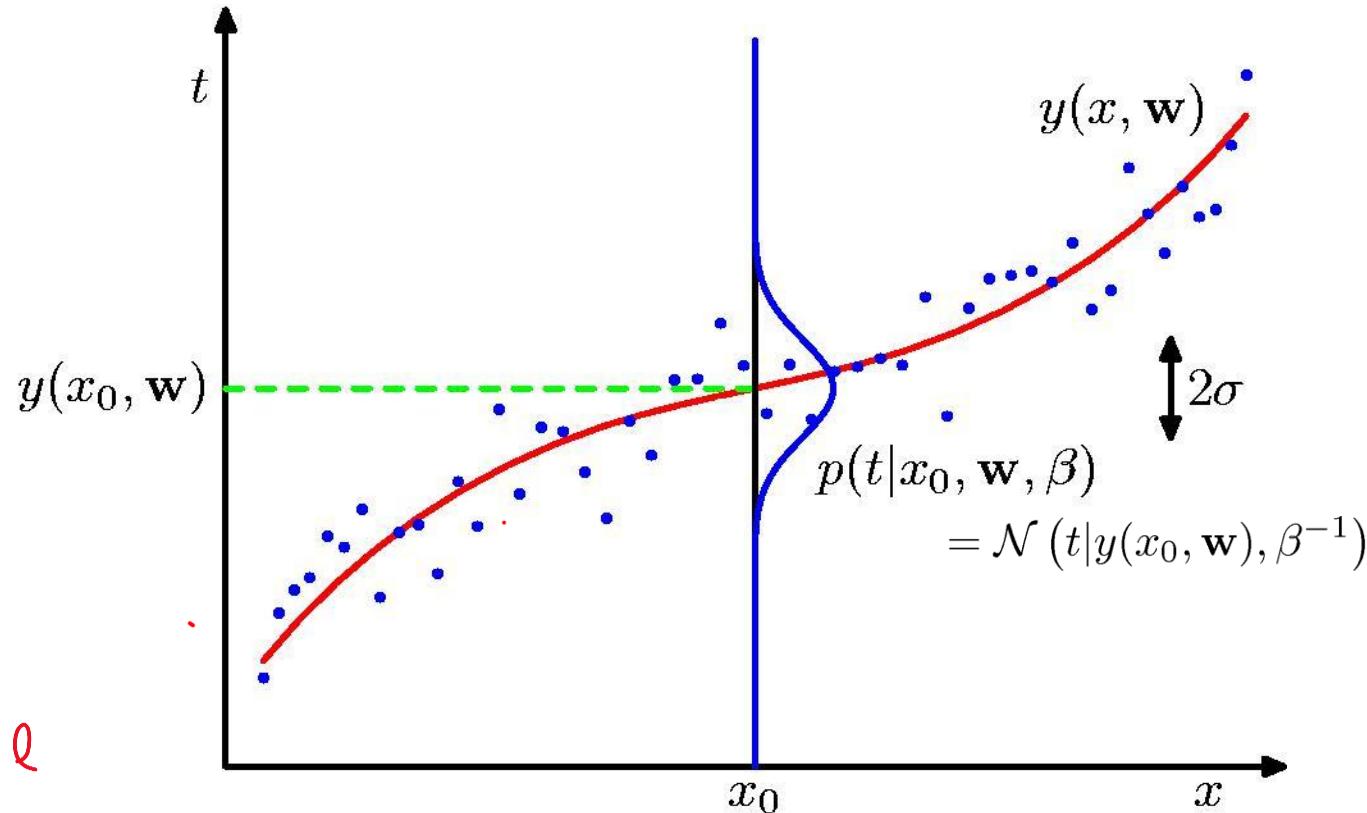
(how? *Density Estimation*; MLE/Bayesian Inference)

(2) Predict t for a new x from estimated $P(t|x)$

(how? *Decision Theory*; e.g., $y(x_{new}) = E[t|x = x_{new}]$)

Curve Fitting: Going to the basics!

using a Probabilistic Model & its Density Estimation (MLE/Bayesian)



*Discriminative
model $t|x$: $t = y(x, \mathbf{w}) + \epsilon$, $\epsilon \sim N(0, \bar{\beta}^2)$, $\bar{\beta} = \sigma^2$*
[CMB]

ML estimation

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

$$\begin{aligned}\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\ &= -\underbrace{\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2}_{\beta E(\mathbf{w})} + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})\end{aligned}$$

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

Determine \mathbf{w}_{ML} by minimizing sum-of-squares error $E_D(w)$.

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

[CMB]

Summary: Linear model for regression -

$$w_{ML} == w_{LS}$$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

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

| S |

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

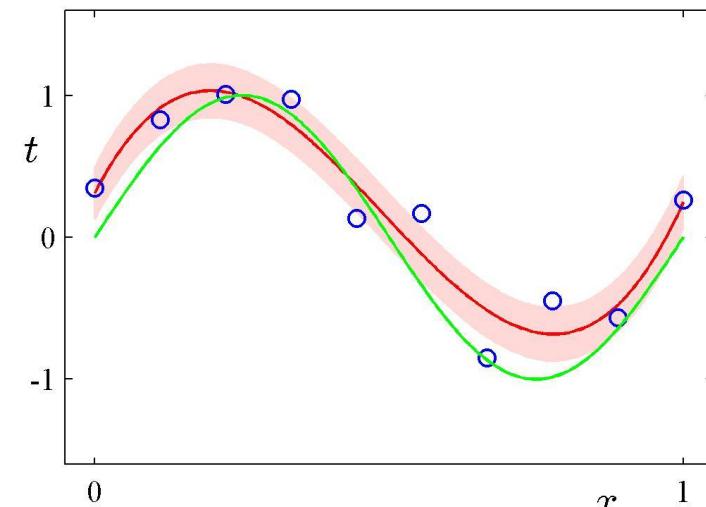
$\Phi^T \Phi \in N \times N$

$\Phi^T \mathbf{t} \in N \times M$

$$t = y(x, \mathbf{w}) + \epsilon$$

$m | \epsilon$

$$p(t|x, \mathbf{w}_{ML}, \beta_{ML}) = \mathcal{N}(t | y(x, \mathbf{w}_{ML}), \beta_{ML}^{-1})$$

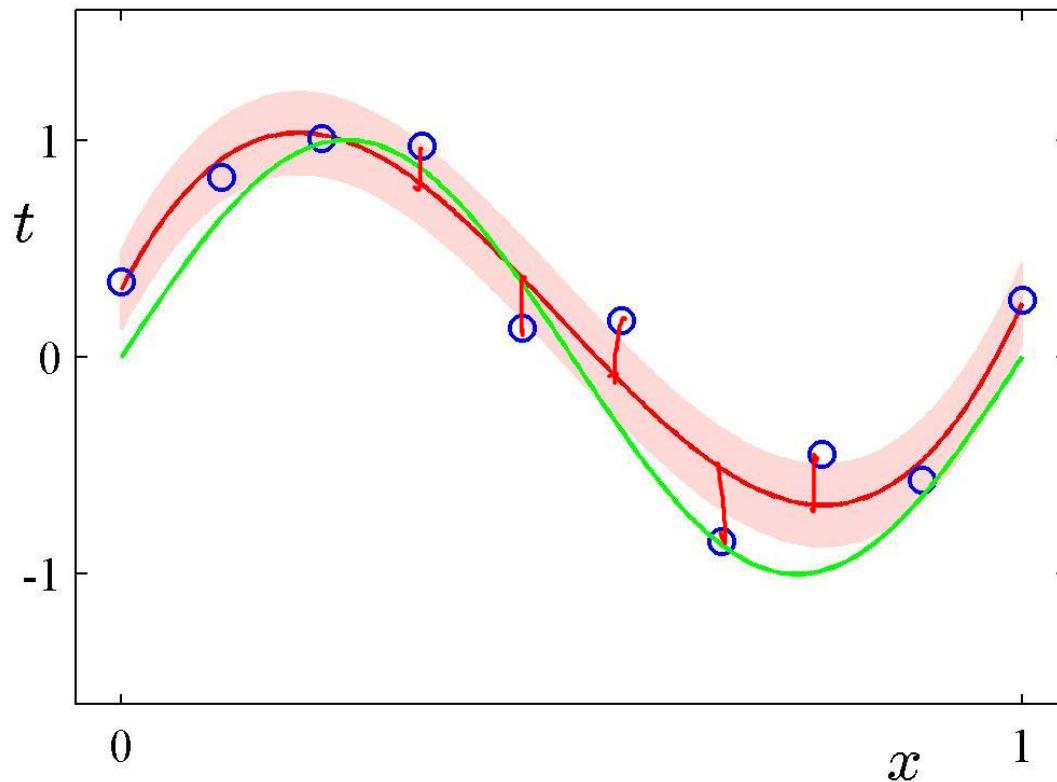


$$\epsilon \sim \mathcal{N}(0, \frac{1}{\beta})$$

[CMB]

Addtnl. Advantage: ML Predictive Distribution

$$p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t|y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1})$$

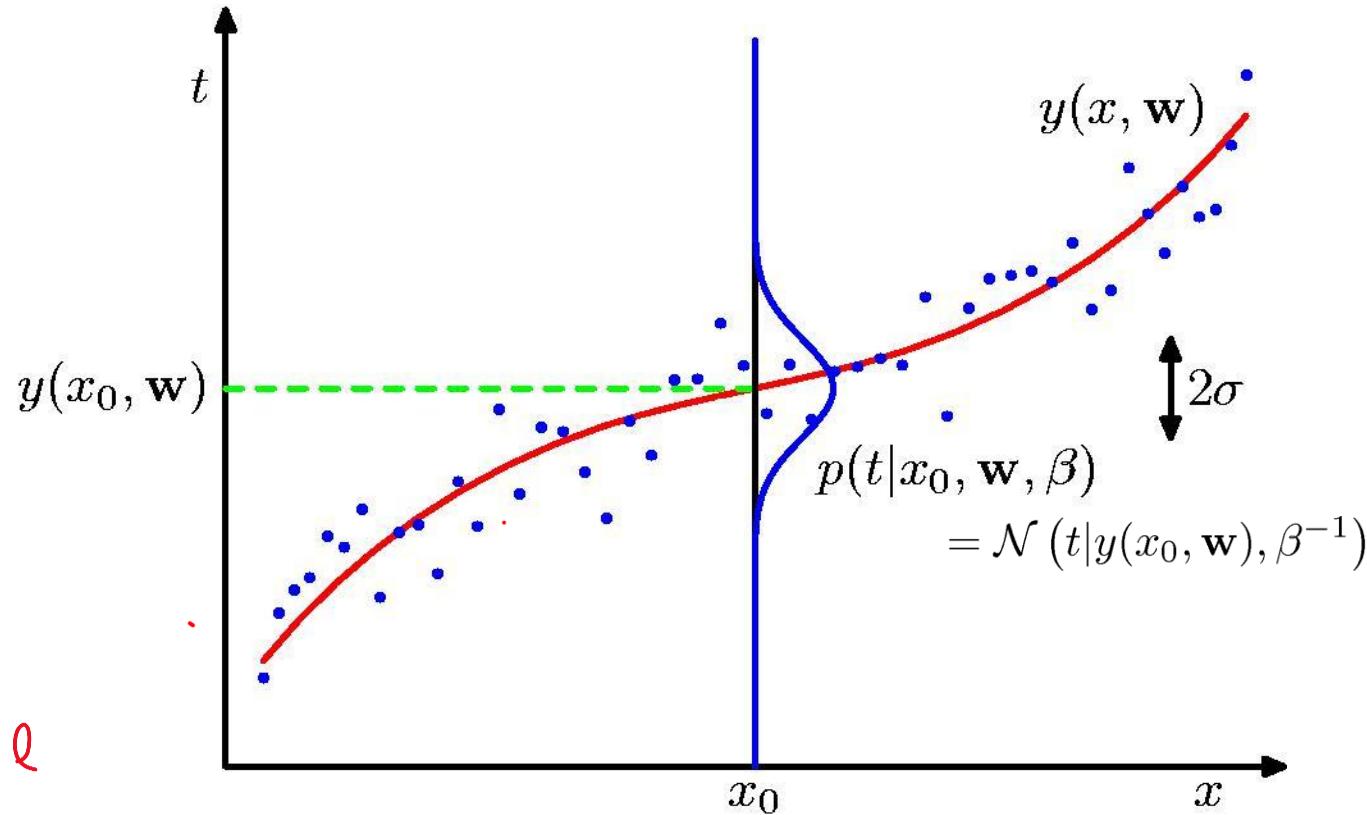


[CMB]

Outline for Module M5

- M6. Regression (Linear models)
 - M6.0 Introduction/context
 - Problem formulation
 - Motivating examples, incl. basis functions
 - **M6.1 Linear regression approaches**
 - Direct approach I: least-squares (error) (w_{LS})
 - Direct approach II: least-squares (error) with regularization (w_{RLS})
 - Discriminative model approach I: MLE ($w_{ML} = w_{LS}$)
 - **Discriminative model approach IIa: Bayesian linear regression ($w_{MAP} = w_{RLS}$)**
 - Discriminative model approach IIb: Bayesian linear regression (posterior & predictions)
 - M6.2 Model Complexity/Selection
 - Motivation (hyperparameter tuning to avoid overfitting)
 - Frequentist view (bias-variance decomposition)
 - M6.3 Concluding thoughts

Recall: Probab. Model & its Density Estimation (MLE)



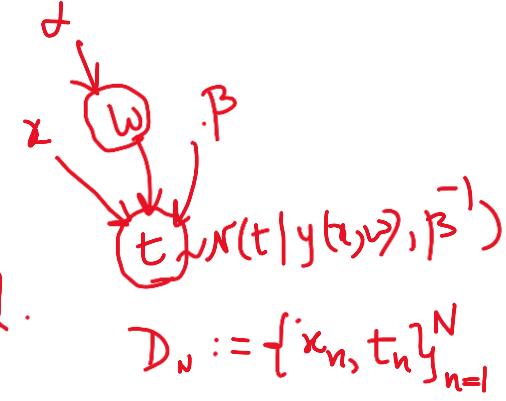
Discriminative
model $t|x$: $t = y(x, \mathbf{w}) + \epsilon$, $\epsilon \sim N(0, \bar{\beta}^{-1})$, $\bar{\beta} = \sigma^2$
[CMB]

Bayesian inference: what would you model as a rv instead of a fixed value?

- Brainstorm
 - What would you model?
 - What would you infer?

Bayesian inference: what would you model as a rv instead of a fixed value?

- Brainstorm
 - What would you model? $P(w)$ prior, besides the $P(t|x, w)$ model.
 - What would you infer?
(1) $P(w|D_N)$ = post.
(2a) MAP-based predn. (or) (2b) $P(t_{\text{new}}|x_{\text{new}}, D_N)$ post. pred.



Compare with MLE:
(1) \hat{w} is same: D_N

0 If desired:
(1) w^* that max. $L(w)$ or min $E(w)$
(2) Pred. $t_{\text{new}} = y(x_{\text{new}}, w^*)$

Relation between Bayesian MAP and Regularized linear regression

Let look at a simpler problem first – mode of the posterior of w

$$w_{MAP} = \operatorname{argmax}_w P(w \mid \mathcal{D}_N := \{\mathbf{x}, \mathbf{t}\} := \{\mathbf{x}_n, \mathbf{t}_n\}_{n=1 \text{ to } N})$$

We will actually show that $w_{MAP} = w_{RLS}$!!

Bayesian inference: a first step via MAP

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

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

Determine maximizer of this posterior, i.e., \mathbf{w}_{MAP} , by minimizing regularized sum-of-squares error $\tilde{E}(\mathbf{w})$, because:

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

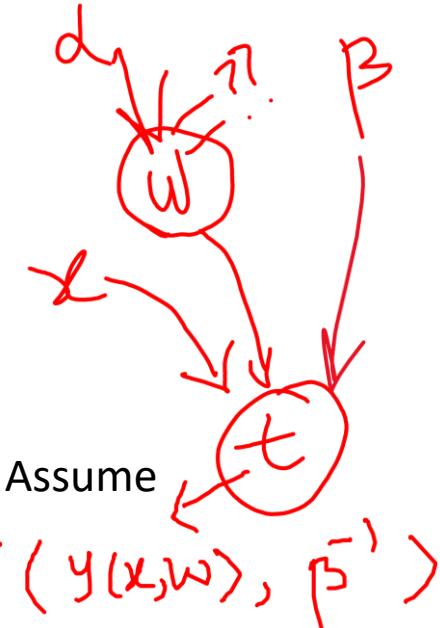
$$\beta E_D(\mathbf{w})$$

contains: likelihood

$$\lambda E_W(\mathbf{w}) \quad \lambda \triangleq \frac{\alpha}{\beta} \cdot \left(\frac{\lambda}{2}\|\mathbf{w}\|^2\right)$$

penalty / regularization / prior
(shrinks \mathbf{w} towards 0)

[cf. full details of proof in next slide]



[CMB]

Full details of w_{MAP} derivation, & $w_{MAP} = w_{RLS}$ proof!

Likhd.

$$p(\underline{t} | \underline{x}, \underline{w}, \beta) = \prod_{n=1}^N p(t_n | x_n, w, \beta)$$

$$= \prod_{n=1}^N \left(\frac{\beta}{2\pi} \right)^{1/2} \exp \left\{ -\frac{\beta}{2} (y(x_n, w) - t)^2 \right\}$$

$$\Rightarrow \ln p(\underline{t} | \underline{x}, \underline{w}, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) \boxed{-\beta E_D(w)}$$

Prior.

$$p(\underline{w} | \alpha) = \left(\frac{\alpha}{2\pi} \right)^{(M+1)/2} \exp \left\{ -\frac{\alpha}{2} \underline{w}^T \underline{w} \right\}$$

$$\ln p(\underline{w} | \alpha) = \frac{M}{2} \ln \alpha - \frac{M}{2} \ln(2\pi) \boxed{-\frac{\alpha}{2} \underline{w}^T \underline{w}}$$

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

$$\begin{aligned} \Rightarrow \ln p(\underline{w} | \underline{x}, \underline{t}, \alpha, \beta) &= \ln p(\underline{t} | \underline{x}, \underline{w}, \alpha, \beta) + \ln p(\underline{w} | \alpha) + \text{const.} && (\text{assume } \alpha, \beta \text{ are both known hyperparams.}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(w) + \frac{M}{2} \ln \alpha - \frac{M}{2} \ln(2\pi) - \frac{\alpha}{2} w^T w + \text{const.} \end{aligned}$$

So maximizing this $\ln(\text{posterior})$ is the same as max. $\boxed{-\beta E_D(w) - \frac{\alpha}{2} \|w\|^2}$ (ignoring terms indept. of w), or equivalently

$$\min. \beta E_D(w) + \alpha \frac{1}{2} \|w\|^2 = \beta \left(E_D(w) + \frac{\alpha}{\beta} E_W(w) \right) = \boxed{\beta \tilde{E}(w)}, \text{ or equivalently min. } \tilde{E}(w). \quad ((\text{here, we set } \lambda := \frac{\alpha}{\beta}))$$

Changing the prior from Gaussian to Laplacian!

- Prior: $p(w) = p(w_0)p(w_1) \dots p(w_{M-1})$
 - What if $p(w_i)$ changed from
$$\left(\frac{\alpha}{2\pi}\right)^{\frac{1}{2}} \exp\left\{-\frac{\alpha}{2} w_i^2\right\} \rightarrow \left(\frac{\alpha}{4}\right) \exp\left\{-\frac{\alpha}{2} |w_i|\right\}$$
?
 - Then, $p(w)$ changes from
$$\left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \exp\left(-\frac{\alpha}{2} \|w\|_2^2\right) \rightarrow \left(\frac{\alpha}{4}\right)^M \exp\left(-\frac{\alpha}{2} \|w\|_1\right)$$

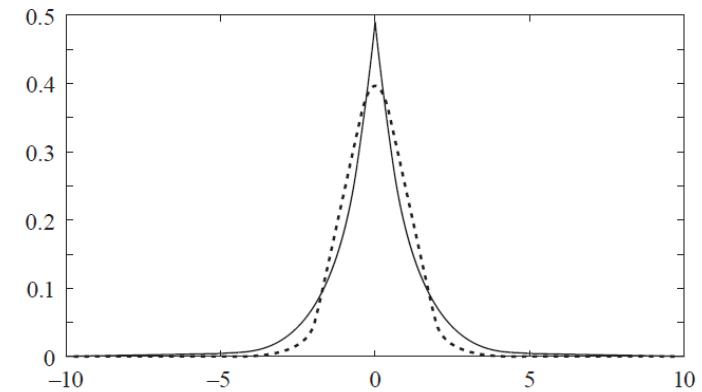


Figure 20.3 Laplacian distribution ($\beta = 1$) and Gaussian distribution ($\sigma^2 = 1$)

- Then regularization term $E_w(w)$ in regularized error $\tilde{E}(w) = E_D(w) + \lambda E_w(w)$ changes from:
 - $\frac{1}{2} \sum_i w_i^2$ (ridge regn. / L_2 regularization) $\rightarrow \frac{1}{2} \sum_i |w_i|$ (lasso regn. / L_1 regularization)

Recall: Our original ridge regn. obj. fn: $\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$

General prior $p(\mathbf{w})$ that generalizes both Gaussian and Laplacian prior

$$p(\mathbf{w}|\alpha) = \left[\frac{q}{2} \left(\frac{\alpha}{2} \right)^{1/q} \frac{1}{\Gamma(1/q)} \right]^M \exp \left(-\frac{\alpha}{2} \sum_{j=1}^M |w_j|^q \right)$$

Outline for Module M5

- M6. Regression (Linear models)
 - M6.0 Introduction/context
 - Problem formulation
 - Motivating examples, incl. basis functions
 - **M6.1 Linear regression approaches**
 - Direct approach I: least-squares (error) (w_{LS})
 - Direct approach II: least-squares (error) with regularization (w_{RLS})
 - Discriminative model approach I: MLE ($w_{ML} = w_{LS}$)
 - Discriminative model approach IIa: Bayesian linear regression ($w_{MAP} = w_{RLS}$)
 - **Discriminative model approach IIb: Bayesian linear regression (posterior & predictions)**
 - M6.2 Model Complexity/Selection
 - Motivation (hyperparameter tuning to avoid overfitting)
 - Frequentist view (bias-variance decomposition)
 - M6.3 Concluding thoughts

Bayesian inference: second and third steps

- Assume Gaussian prior for w going forward.
- We don't want just a single-point estimate (MAP) of w ; we want

Step 2) the full posterior of w , and in turn use it to

Step 3) predict t for new x (via model-averaging...

- ...wherein each model is given by a particular possible value of w and the averaging weight is given by the model's posterior)

Step 2) Let's see an example of full posterior...

$$t \cong y(x, w) = w_0 + w_1 x$$

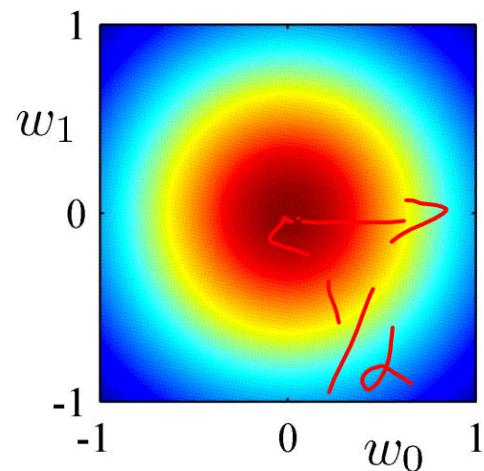
Use a given set of training data points to not just infer one optimal model specified by w_{MAP} , but all possible w models and the training dataset's support (posterior probab.) for each such model w .

Bayesian Linear Regression Example (1)

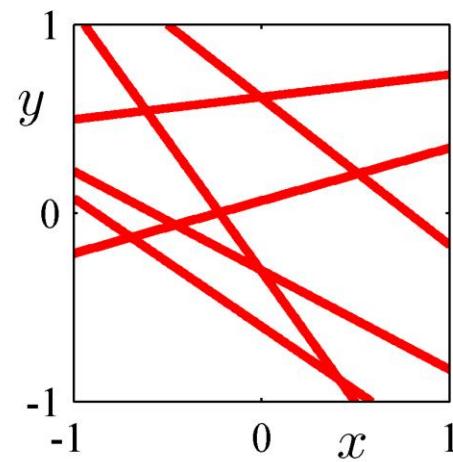
0 data points observed

$$N(0, \sigma^2 I)$$

Prior

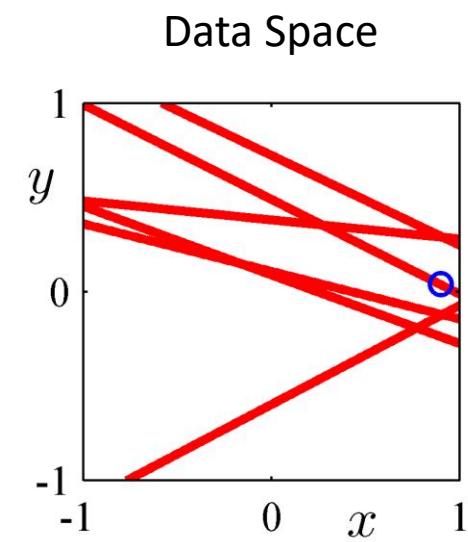
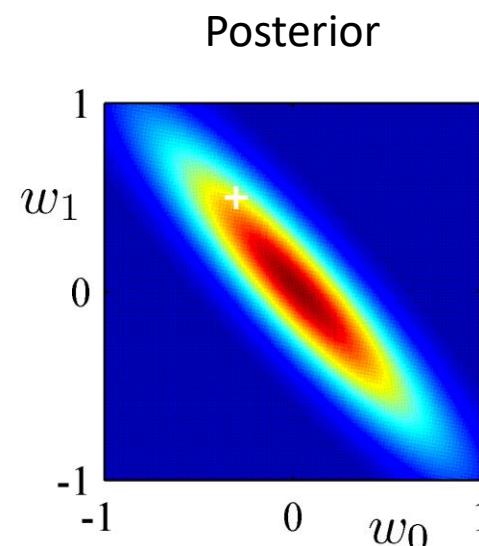
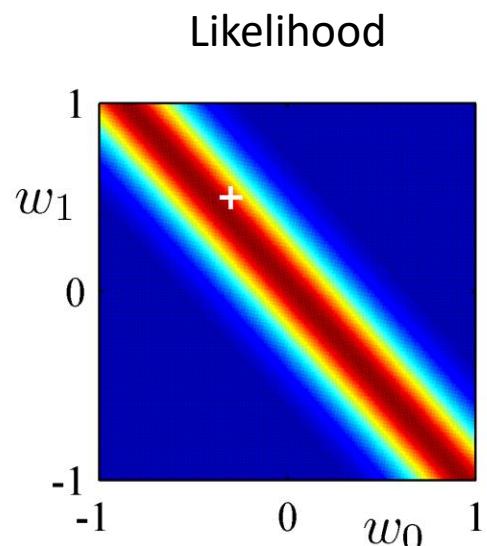


Data Space



Bayesian Linear Regression Example (2)

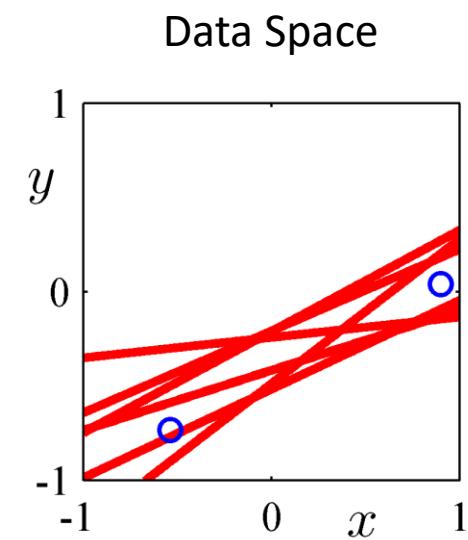
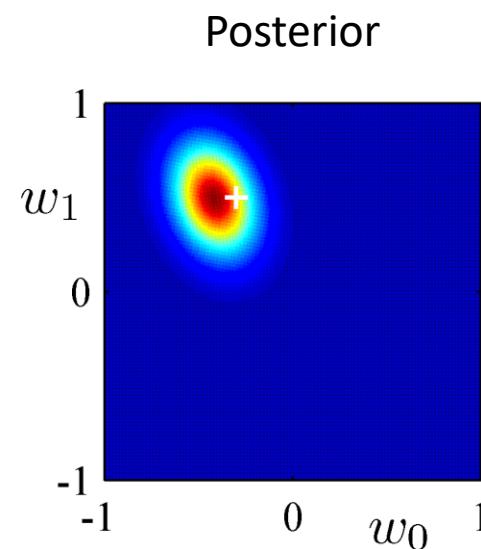
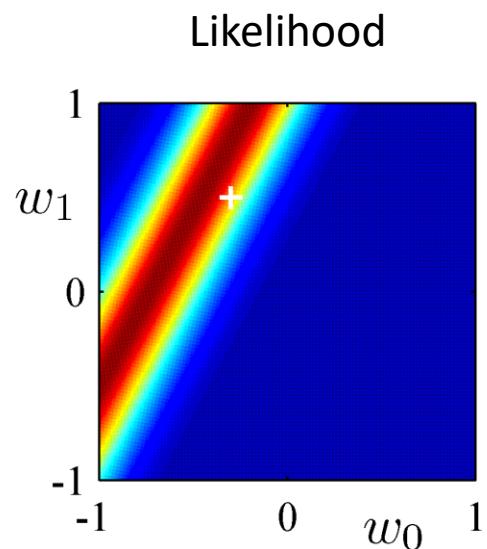
1 data point observed



$$P(w|t) \propto P(t|w) P(w)$$

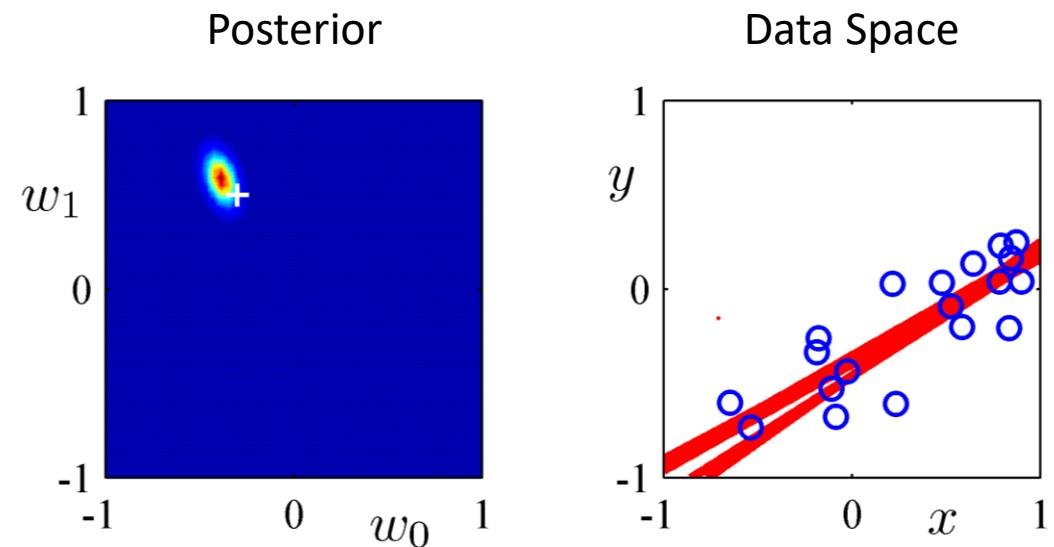
Bayesian Linear Regression Example (3)

2 data points observed



Bayesian Linear Regression Example (4)

20 data points observed



From example posterior plots to
Full posterior in the general case:

Bayesian Linear Regression

- Define a conjugate **prior** over w . A common choice is

$$p(w) = \mathcal{N}(w | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

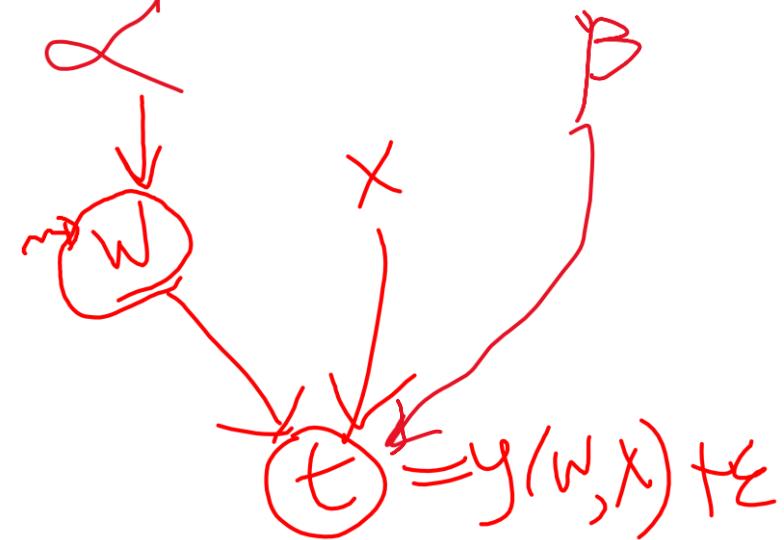
- Combining this with the **likelihood** function and using results for marginal and conditional Gaussian distributions, gives the **posterior**

$$p(w | \mathbf{D}) = \mathcal{N}(w | \mathbf{m}_N, \mathbf{S}_N)$$

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi.$$

[Qn. What is m_N ?]



Marginal and Conditional Gaussians

Given a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$\begin{aligned} p(x) &= \mathcal{N}(x | \mu, \Lambda^{-1}) \\ p(y|x) &= \mathcal{N}(y | Ax + b, L^{-1}) \end{aligned} \quad (2.113) \quad (2.114)$$

the marginal distribution of y and the conditional distribution of x given y are given by

$$\begin{aligned} p(y) &= \mathcal{N}(y | A\mu + b, L^{-1} + A\Lambda^{-1}A^T) \\ p(x|y) &= \mathcal{N}(x | \Sigma \{ A^T L (y - b) + A\mu \}, \Sigma) \end{aligned} \quad (2.115) \quad (2.116)$$

where

$$\Sigma = (\Lambda + A^T L A)^{-1}. \quad (2.117)$$

[CMB]

Recall: MVG Handy Results (cheat-sheet)



Partitioned Gaussians

Given a joint Gaussian distribution $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$ and

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad (2.94)$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}, \quad \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}. \quad (2.95)$$

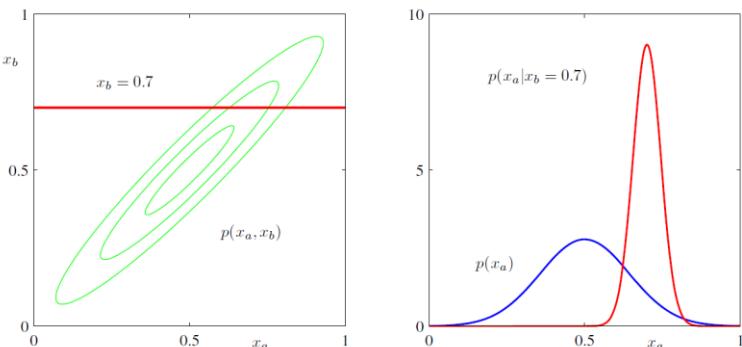
Conditional distribution:

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1}) \quad (2.96)$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b). \quad (2.97)$$

Marginal distribution:

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}). \quad (2.98)$$



Marginal and Conditional Gaussians

Given a marginal Gaussian distribution for \mathbf{x} and a conditional Gaussian distribution for \mathbf{y} given \mathbf{x} in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (2.113)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|A\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \quad (2.114)$$

the marginal distribution of \mathbf{y} and the conditional distribution of \mathbf{x} given \mathbf{y} are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|A\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \quad (2.115)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \quad (2.116)$$

where

$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1}. \quad (2.117)$$

Bang Lin Legn.

[CMB: Bishop, Chapter 2]

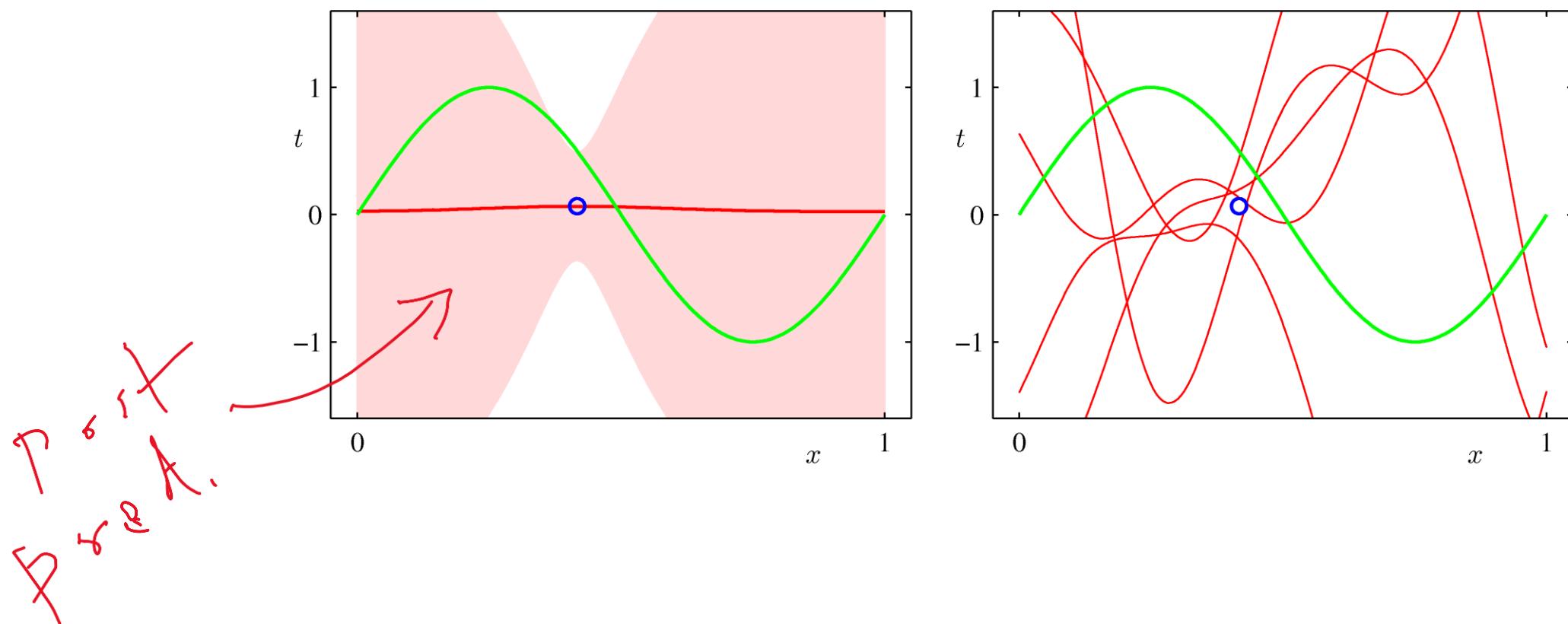
Finally, Step 3: What about predicting t for a new point x ?

- ((we don't want just an estimate or posterior of w ; we want to use it to predict t for new x))

((again example plots first, and post. predictive distbn. form for the general case after that.))

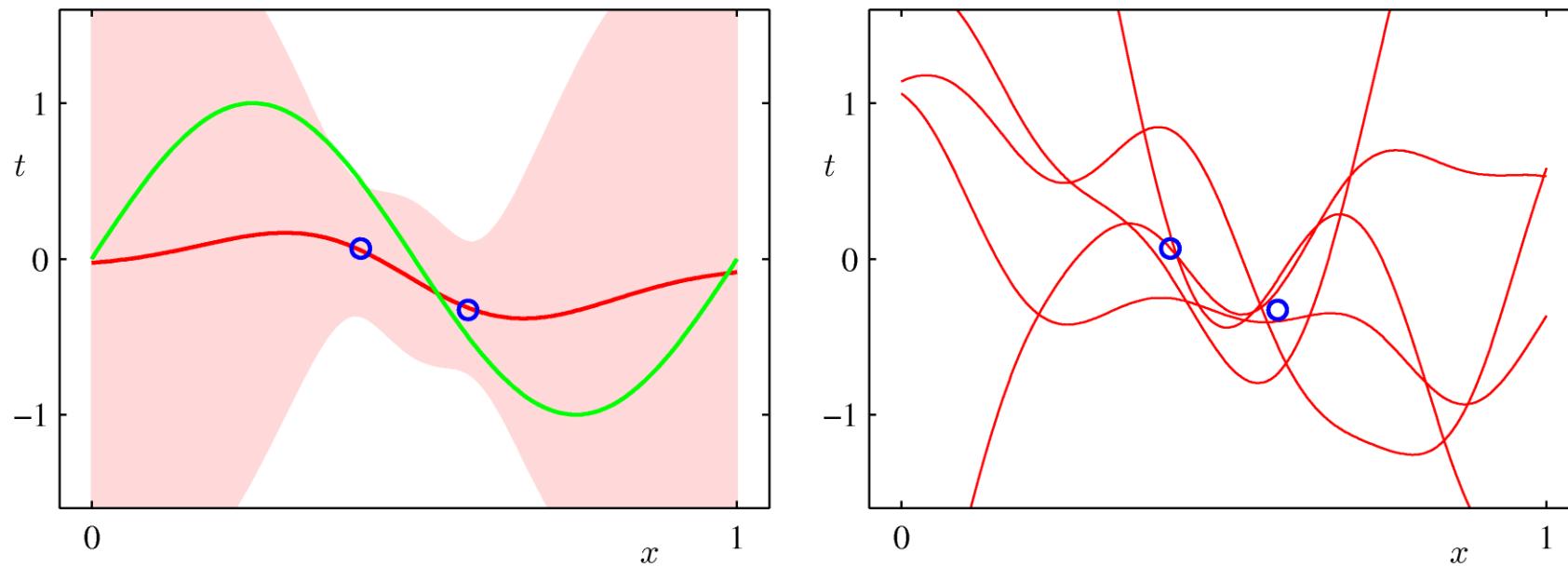
Example Predictive Distribution (1)

- Example: Sinusoidal data, 9 Gaussian basis fns. ($\phi: \mathbb{R} \rightarrow \mathbb{R}^9$), 1 data point



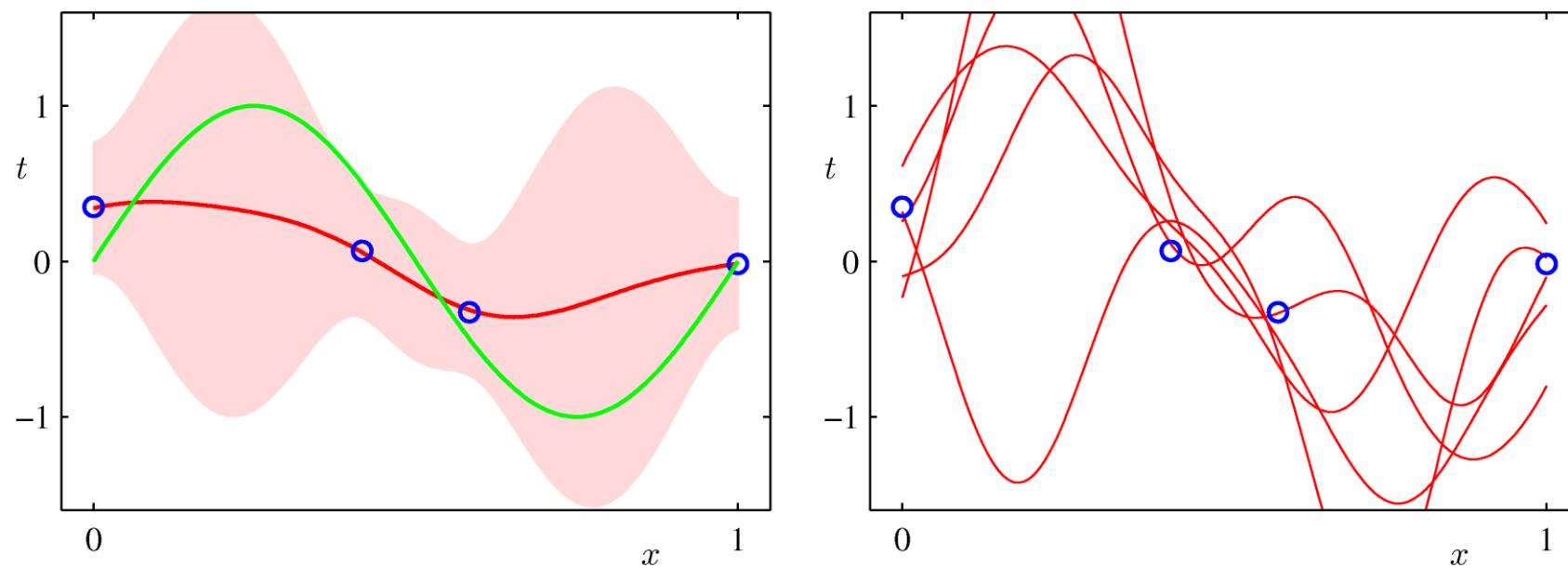
Example Predictive Distribution (2)

- Example: Sinusoidal data, 9 Gaussian basis functions, 2 data points



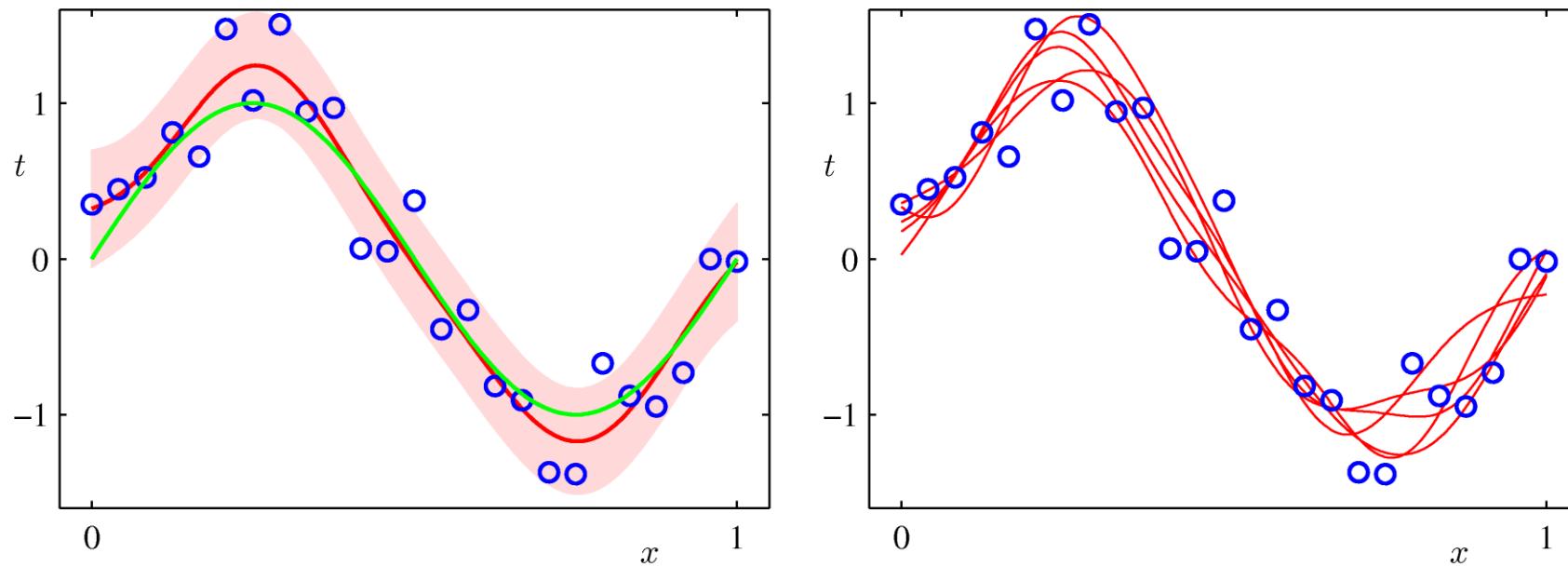
Example Predictive Distribution (3)

- Example: Sinusoidal data, 9 Gaussian basis functions, 4 data points



Example Predictive Distribution (4)

- Example: Sinusoidal data, 9 Gaussian basis functions, 25 data points



(Bayesian/Posterior) Predictive Distribution (1)

- Predict t for new values of x by integrating over w (model-averaging)

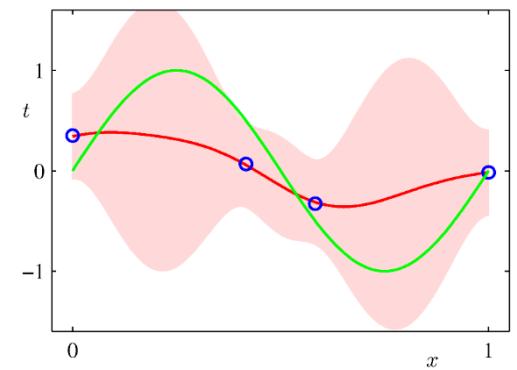
$$\begin{aligned}
 p(t | x, \underline{\mathbf{x}}, \underline{\mathbf{t}}, \alpha, \beta) &= \int_{\omega} p(t | x, \mathbf{w}, \beta) \underset{\text{post.}}{p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta)} d\mathbf{w} \\
 &= \mathcal{N}(t | \mathbf{m}_N^T \boldsymbol{\phi}(x), \underset{\text{uncert.}}{\underline{\sigma_N^2(x)}})
 \end{aligned}$$

- where

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \underbrace{\boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x)}_{(\text{ML})} \underset{\text{(extra due to uncertainty in } \omega\text{)}}{\text{extra due to uncertainty in } \omega}$$

$$\begin{aligned}
 \mathbf{m}_N &= \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} \\
 \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}.
 \end{aligned}$$

Exercise: Prove that $\mathbf{m}_N = w_{RLS} = w_{MAP}$ (recall $\lambda := \alpha/\beta$), and hence that the posterior predictive mean is same as that of the predicted value in the direct RLS approach. [CMB]



Linear regression direct vs. discriminative model approaches: summary

Discriminative model-based approaches have two advantages:

1. Convert intuition for obj. fns. to probabilistic model driven motivations:

$$(\text{Least-squares or min } E(w)) \quad w_{LS} = w_{ML} \quad (\text{MLE})$$

$$(\text{Reg. Least-squares or min } \tilde{E}(w)) \quad w_{RLS} = w_{MAP} \quad (\text{Bayesian})$$

2. Give additional advantage of capturing the uncertainty over the predicted values, viz., a predictive **distribution**

$p(t|x) = N(t | y(x, w_*) := w_*^T \phi(x), \mathbf{var})$, (and **not just a single predicted value** $y(x, w_*)$ as in the direct approach).

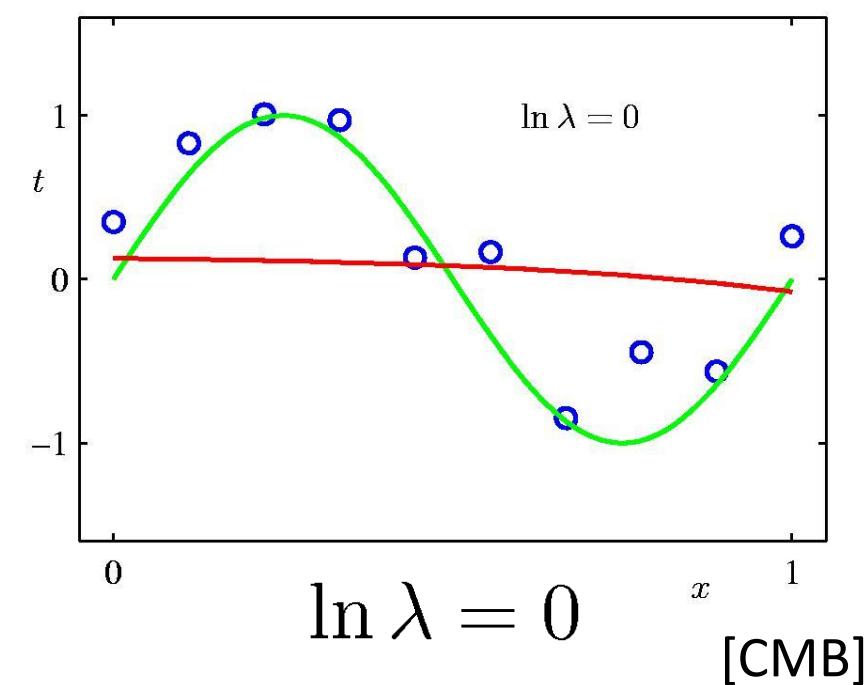
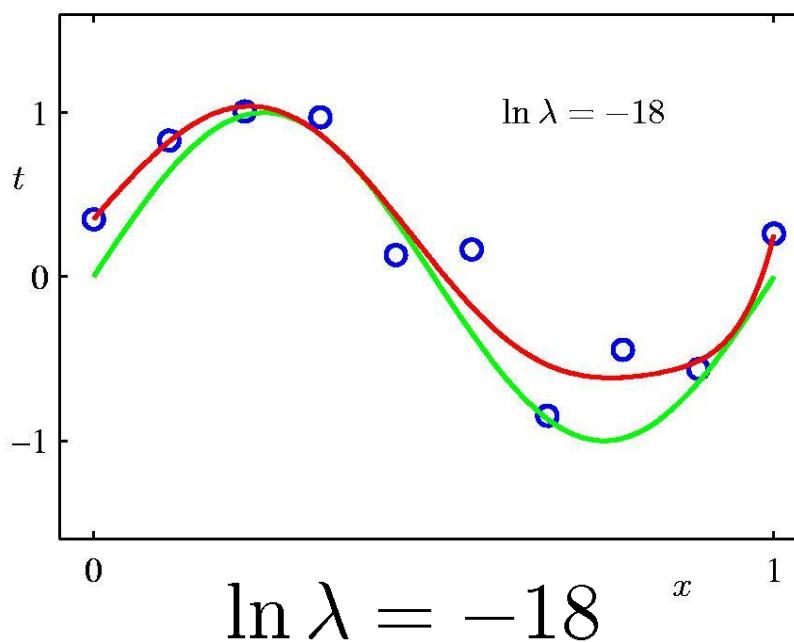
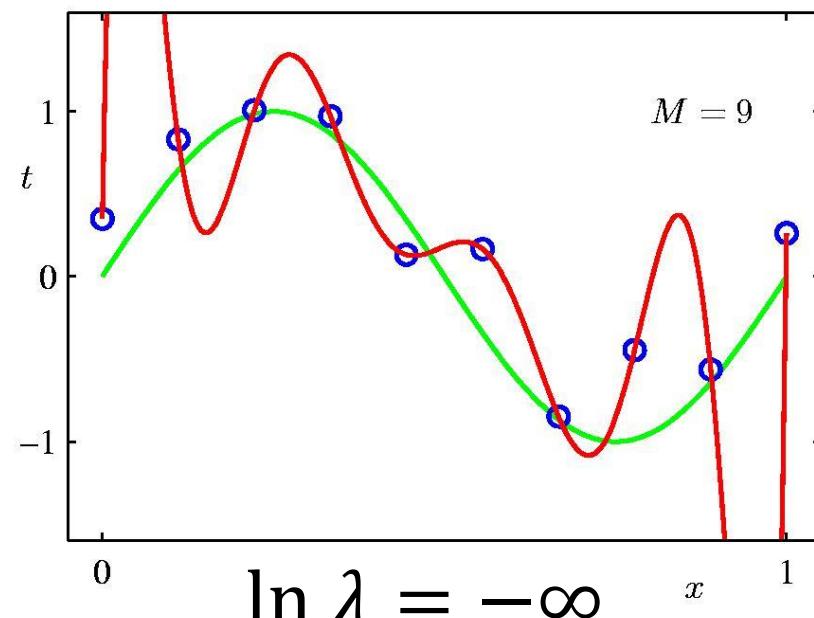
- In MLE, (pred.) **var** is a dataset-wide single variance ($\sigma^2 = \beta^{-1}$)
- In Bayesian, (post. pred.) **var** is datapoint-specific ($\sigma^2(x) = \beta^{-1} + \phi(x)^T S_N \phi(x)$)

Outline for Module M5

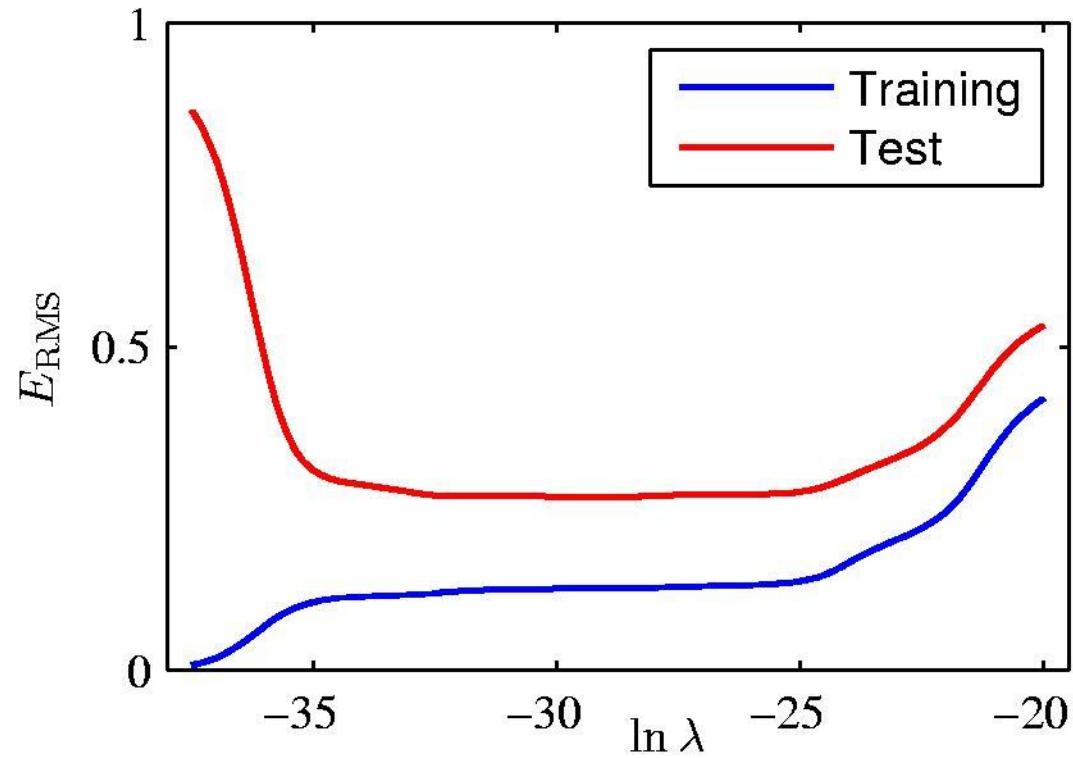
- M6. Regression (Linear models)
 - M6.0 Introduction/context
 - Problem formulation
 - Motivating examples, incl. basis functions
 - M6.1 Linear regression approaches
 - Direct approach I: least-squares (error) (w_{LS})
 - Direct approach II: least-squares (error) with regularization (w_{RLS})
 - Discriminative model approach I: MLE ($w_{ML} = w_{LS}$)
 - Discriminative model approach IIa: Bayesian linear regression ($w_{MAP} = w_{RLS}$)
 - Discriminative model approach IIb: Bayesian linear regression (posterior & predictions)
 - **M6.2 Model Complexity/Selection**
 - **Motivation (hyperparameter tuning to avoid overfitting)**
 - Frequentist view (bias-variance decomposition)
 - M6.3 Concluding thoughts

Recall: Motivating example for Regularization, and Hyperparam. Tuning

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



Regularization, and Hyperparam. tuning



Regularization penalizes large coeffs.
of 9th order polynomial fit of the data.

more complex model ← [CMB]

less complex model

Some Motivating Questions

- Can we understand the error in our predictions better?
 - That is, can we identify the different components of the error in our predictions?
 - How are these different components related to the complexity of our model, and to overfitting?
- Can we use above knowledge to better tune model complexity (hyperparameters) to avoid overfitting?
- If the trends in data require fitting of a complex model, then
 - can overfitting be detected by understanding the stability of the optimal (frequentist/ML) model across different training datasets?
 - can a Bayesian model overcome the overfitting “naturally/implicitly” by not settling in on a single optimal model and instead averaging over multiple models?
 - Bayesian view (model averaging & empirical Bayes) possible, but out of scope for this course.

Outline for Module M5

- M6. Regression (Linear models)
 - M6.0 Introduction/context
 - Problem formulation
 - Motivating examples, incl. basis functions
 - M6.1 Linear regression approaches
 - Direct approach I: least-squares (error) (w_{LS})
 - Direct approach II: least-squares (error) with regularization (w_{RLS})
 - Discriminative model approach I: MLE ($w_{ML} = w_{LS}$)
 - Discriminative model approach IIa: Bayesian linear regression ($w_{MAP} = w_{RLS}$)
 - Discriminative model approach IIb: Bayesian linear regression (posterior & predictions)
 - **M6.2 Model Complexity/Selection**
 - Motivation (hyperparameter tuning to avoid overfitting)
 - **Frequentist view (bias-variance decomposition)**
 - **M6.3 Concluding thoughts**

Recall: Decision Theory for Regression
min. squared loss (cond. expn. as minimizer)

$$\mathbb{E}[L] = \mathbb{E}_x \left[\mathbb{E}_{t|x} [(y(x) - t)^2] \right]$$

$$\mathbb{E}[L] = \underbrace{\int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x}}_{(y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}))^2} + \underbrace{\int \text{var}[t|\mathbf{x}] p(\mathbf{x}) d\mathbf{x}}_{\text{noise}}$$

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt.$$

Bias-variance analysis proof

$$\mathcal{D} := \mathcal{D}_N = \left\{ \begin{bmatrix} x_n \\ t_n \end{bmatrix} \right\}_{n=1}^N$$

$$\left(\begin{bmatrix} x_n \\ t_n \end{bmatrix} \right) \stackrel{iid}{\sim} P(x, t)$$

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt.$$

$$\mathbb{E}[L] = \int \underbrace{\{y(\mathbf{x}) - h(\mathbf{x})\}^2}_{(y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}))^2} p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt.$$

Goal: Decompose average error $E_{\mathcal{D}}[E_{x,t}[L]]$ into different terms.

Now, simply view “ $y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})$ ” as a random variable Z ; and apply the variance formula:
 $Var_{\mathcal{D}}(Z) = E_{\mathcal{D}}[Z^2] - (E_{\mathcal{D}}[Z])^2$ to get the bias-variance decomposition of error below:

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ &= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}. \end{aligned}$$

Bias-variance decomposition in formula

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

where

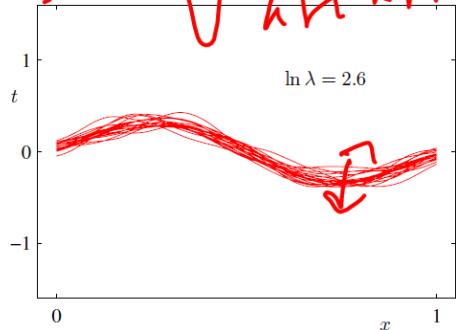
$$\begin{aligned} (\text{bias})^2 &= \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} & h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt. \\ \text{variance} &= \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x} \\ \text{noise} &= \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \end{aligned}$$

$$\text{expected loss} = E_{\mathcal{D}} [E_{x,t}[L]] = E_{\mathcal{D}} [E_{x,t}[(y(x; D) - t)^2]]$$

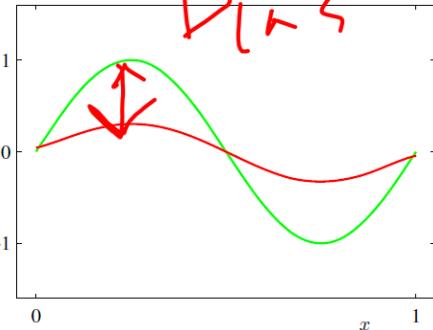
Exercise: cf. worksheet for careful understanding of what random variables the expectation above is taken over!

Bias-variance in pictures (for an example)

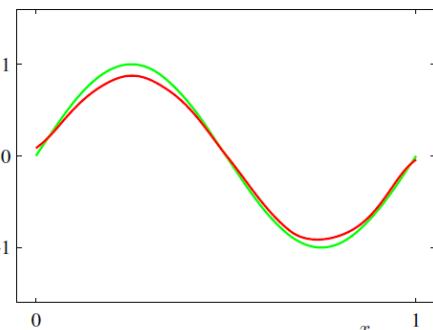
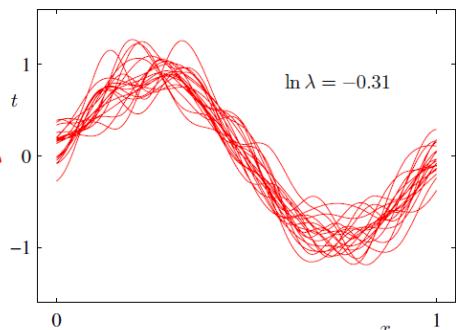
Complexity



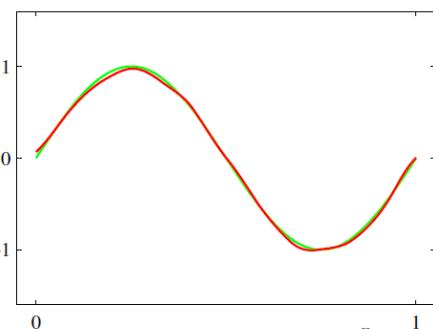
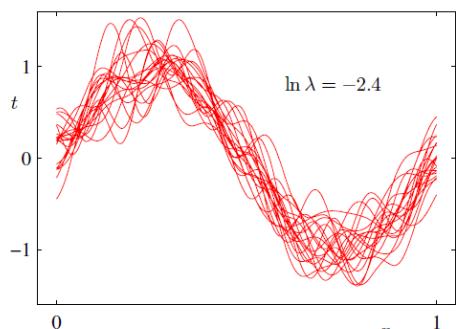
Variance



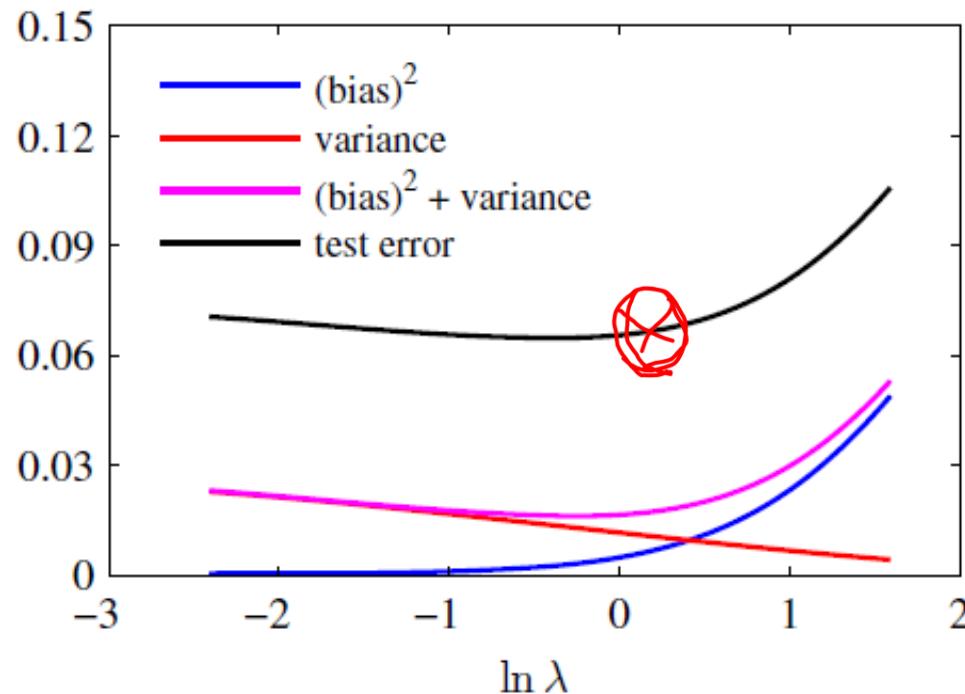
Complexity



Complexity



Bias-variance analysis (for the example)



h_i ←
model opt |
 l_0

expected loss = $(\text{bias})^2 + \text{variance} + \text{noise}$

[CMB]

Bias-variance anal.: applicability in practice?

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

where

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] ??$$

.

[CMB]

Bias-variance anal.: applicability in practice? details

$$\mathbb{E}_{\mathcal{D}}[\mathbb{E}_{x,t}[L]]$$

expected loss = (bias)² + variance + noise

where

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

$$\bar{y}(x) = \frac{1}{L} \sum_{l=1}^L y^{(l)}(x)$$

$$(\text{bias})^2 = \frac{1}{N} \sum_{n=1}^N \{\bar{y}(x_n) - h(x_n)\}^2$$

$$\text{variance} = \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \{y^{(l)}(x_n) - \bar{y}(x_n)\}^2$$

$$h(x) = \mathbb{E}[t|x] ??$$

$\{x_n\}_n^N \sim P(x)$ \mathbb{E}_x \mathbb{E}_t

[CMB]

Concluding thoughts

- Linear regression forms the foundations of other sophisticated methods, so it is good to invest enough time on it.
 - Two views: direct loss fn. view ($E(w)$)/regularization & probabilistic model view (MLE/Bayesian)
 - But lin. regn. has limitations in practice, even with non-linear basis functions, closed-form solutions and other analytical advantages.
 - Mainly because basis functions are fixed before seeing the training data (curse of dimensionality when dimensionality of feature vectors D grows).
- Next steps:
 - linear models for classification, which play similar basic role for other classification methods.
 - Move from fixed basis fns to selection of basis functions or adaptation of basis functions using training data, in later lectures on non-linear models.

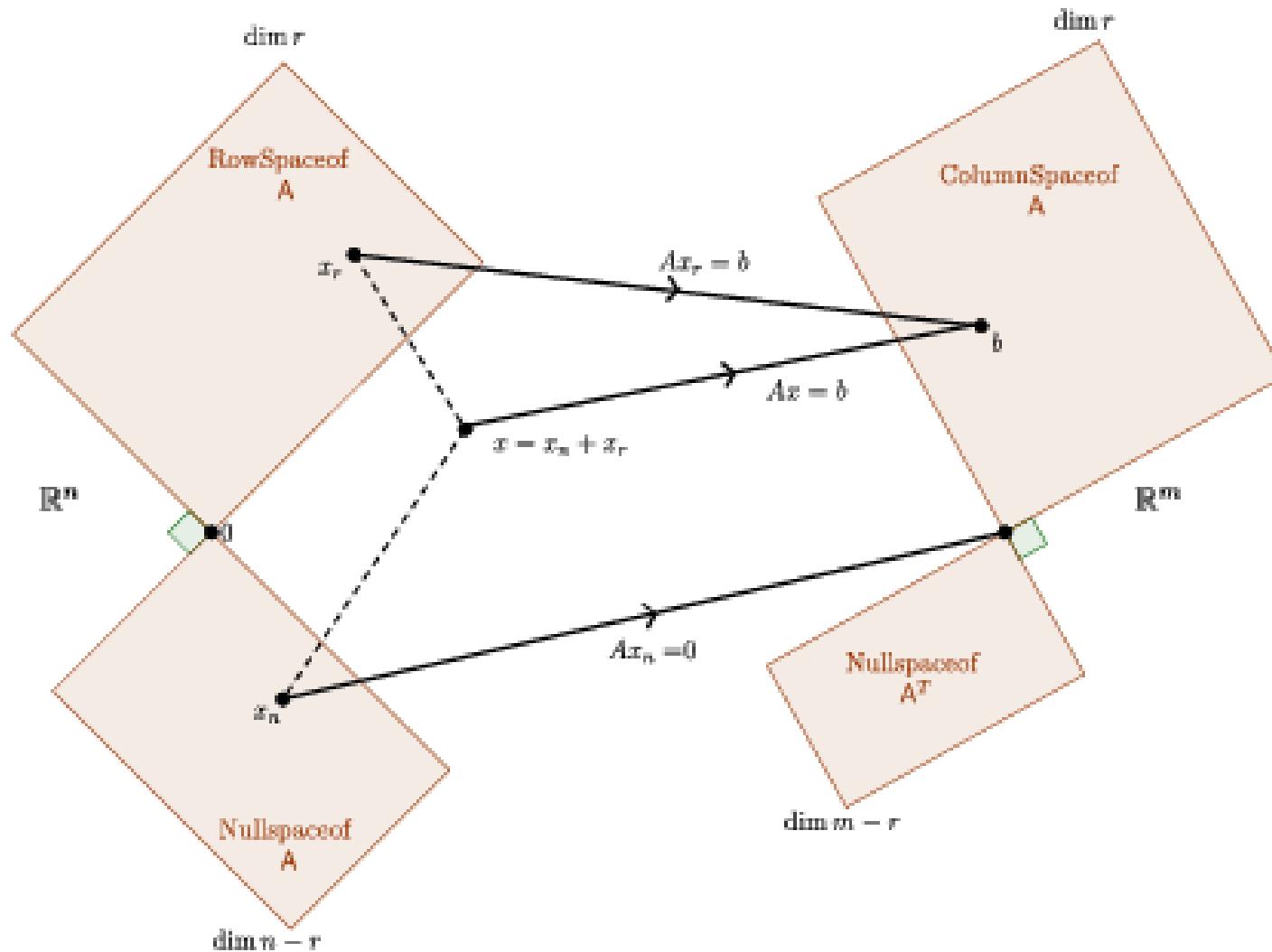
Thank you!

Backup slides

Linear Algebra (LA) Refresher

- Switch to LN Pandey's notes

LA Cheat Sheet: The four subspaces of a matrix



[From <https://mathworld.wolfram.com/FundamentalTheoremofLinearAlgebra.html>, Strang LA book]

LA + Opt. Cheat Sheet

- Real, symmetric matrices S can be diagonalized as $S = Q\Lambda Q^T$
 - S is psd iff all its eigen values are non-negative.

- $$\text{RowSpace}(A) \oplus \text{NullSpace}(A) = \mathbb{R}^d \quad A \in \mathbb{R}^{n \times d}$$

$$\text{ColSpace}(A) \oplus \text{LeftNull}(A) = \mathbb{R}^n$$

$$\therefore \dim(\text{RS}(A)) + \dim(\text{NS}(A)) = d$$

$$\dim(\text{CS}(A)) + \dim(\text{LNS}(A)) = n$$

$$\dim(\text{RS}(A)) = \dim(\text{CS}(A))$$

- - i) $f(x) = x^T Ax \quad \text{for some } A \in \mathbb{R}^{d \times d}$
 $\nabla f(x) = \begin{cases} 2Ax & \text{if } A \text{ is symm} \\ A^T x + Ax & \text{o.w.} \end{cases}$
 - ii) $f(x) = w^T x \quad \text{for some } w \in \mathbb{R}^d$
 $\nabla f(x) = w$
 - iii) $(AB)^T = (B^T A^T)$
 - iv) $(AB)^{-1} = B^{-1} A^{-1} \quad \text{if } AB \text{ is invertible.}$
 - v) If A is not invertible $\Leftrightarrow \exists x \text{ st } Ax = 0$
 $\exists x \quad x^T \cancel{Ax} = 0$
 - vi) $f(x) = g(Ax) \quad f: \mathbb{R}^d \rightarrow \mathbb{R}, g: \mathbb{R}^n \rightarrow \mathbb{R}$
 $\nabla f(x) = A^T \nabla g(Ax) \quad A \in \mathbb{R}^{n \times d}$

[HR]

Recall LA: To solve $Ax = b$, we premult. by A^T , and simply solve $A^T A x = A^T b$.

Linear Equation Solving

$$Ax = b$$

$\textcircled{1} \rightarrow$ No solution if $b \notin \text{column space}(A)$

\rightarrow Unique solution if $b \in \text{column space}(A)$ & A has lin. ind. columns.

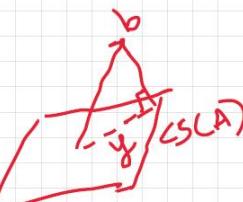
\rightarrow Infinite solutions if $b \in \text{CS}(A)$ & A has lin. dep. columns.

$$\begin{bmatrix} 1 & 3 \\ 2 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$\textcircled{2}$ Least-squares soln possible (related to lin. regn)

$$\min_x \|Ax - b\|^2 = \min_{y \in \text{CS}(A)} \|y - b\|^2$$

$(Ax^* = y^*)$ (y^* is proj. of b onto $\text{CS}(A)$)



$$\min_{x \in \mathbb{R}^n} [(Ax - b)^T (Ax - b)] \rightarrow f(x)$$

$$\nabla f(x) = 2A^T(Ax - b) \stackrel{\text{set to } 0}{=} 0$$

$$\Rightarrow A^T(Ax^* - b) = 0$$

$$\Rightarrow A^T A x^* = A^T b$$

$$\Rightarrow x^* = (A^T A)^{-1} A^T b$$

(normal eqn.)

Ex.: Prove:

- 1) at least one soln. x^* exists for the normal eqn.
- 2) soln. x^* unique if $(A^T A)$ is invertible ($\Leftrightarrow A$ has lin. indep. cols.)
- 3) infinite solns. x^* if $(A^T A)$ is non-invertible ($\Leftrightarrow A$ has lin. dep. cols.)

Ex.:

i) Prove $NS(A) = NS(A^T A)$.

ii) Use orthog. complementarity of $NS(A^T)$, $CS(A)$ to derive normal eqns.

Choice of Prior

Different priors on parameter θ ($\theta := w$) leads to...

$$P(\theta | \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{\theta_i^2}{2\sigma^2} \right\},$$

$$P_{Laplacian}(\theta | \beta) = \frac{1}{2\beta} \exp \left\{ -\frac{|\theta|}{\beta} \right\}.$$

$$P(\theta | \beta) = \prod_{i=1}^n \frac{1}{2\beta} e^{-\frac{|\theta_i|}{\beta}}$$

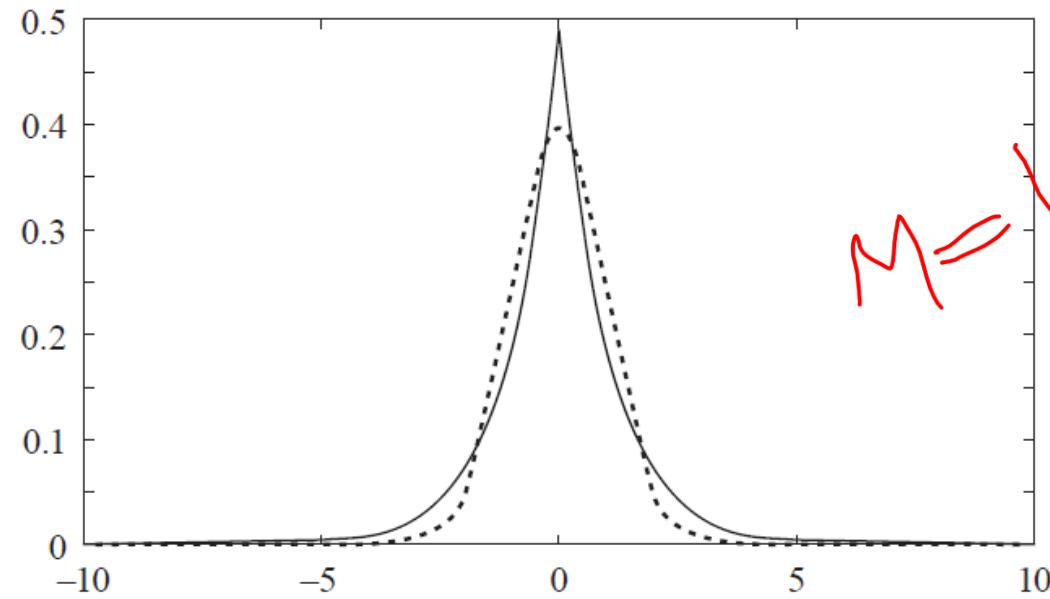


Figure 20.3 Laplacian distribution ($\beta = 1$) and Gaussian distribution ($\sigma^2 = 1$)

...different regularizations (ridge vs. lasso regn.)

$$\text{Ansatz: } \sum_{i=1}^m + \frac{\theta_i}{2\sigma^2} \min_{\text{ridge}}$$

$$|\gamma| \cdot \sum_{i=1}^m + \frac{|\theta_i|}{\beta} \min_{\text{lasso}}$$

Bias-variance analysis (alternate proof)

$$\mathcal{D} := \mathcal{D}_N = \left\{ \mathbf{x}_n, t_n \right\}_{n=1}^N$$

$$(\mathbf{x}_n, t_n) \stackrel{iid}{\sim} P(\mathbf{x}, t)$$

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt.$$

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt.$$

$$\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2.$$

$$\begin{aligned} & \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ & \quad + 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}. \end{aligned}$$

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ &= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}. \end{aligned}$$