

EN.600.475 Machine Learning

Regularization

Raman Arora
Lecture 8
February 20, 2017

- Model complexity and overfitting
- Shrinkage methods: ridge regression, Lasso

Slides credit: Greg Shakhnarovich

Review

Review: noise model and log-likelihood

- Statistical model: noise as a Gaussian random variable

$$y = f(\mathbf{x}; \mathbf{w}) + \nu, \quad \nu \sim \mathcal{N}(\nu; 0, \sigma^2)$$

equivalent to $p(y|\mathbf{x}; \mathbf{w}, \sigma) = \mathcal{N}(y; f(\mathbf{x}; \mathbf{w}), \sigma^2)$

- Maximizing log-likelihood under this model

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} \sum_i \log p(y_i | \mathbf{x}_i; \mathbf{w}, \sigma)$$

is equivalent to least-squares regression

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \sum_i (y_i - f(\mathbf{x}_i; \mathbf{w}))^2$$

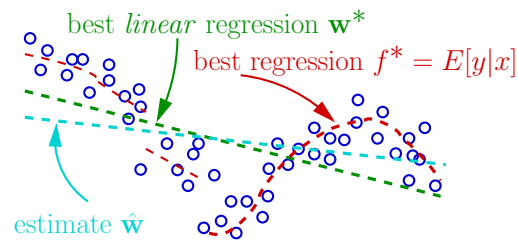
Review: Decomposition of error

- Approximation error

$$E \left[\left(y - \mathbf{w}^{*T} \mathbf{x} \right)^2 \right]$$

- Estimation error

$$E \left[\left(\mathbf{w}^{*T} \mathbf{x} - \hat{\mathbf{w}}^T \mathbf{x} \right)^2 \right]$$



- Approximation error: due to the failure to include optimal predictor in the model class, plus inherent uncertainty in $y|\mathbf{x}$
- Estimation error: due to failure to select the best predictor in the chosen model class; could be reduced with more data

3



Review: generalized linear regression

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 \phi_1(\mathbf{x}) + w_2 \phi_2(\mathbf{x}) + \dots + w_m \phi_m(\mathbf{x}),$$

- Still the same ML estimation technique applies:

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

where \mathbf{X} is the *design matrix*

$$\begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \dots & \phi_m(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \phi_2(\mathbf{x}_2) & \dots & \phi_m(\mathbf{x}_2) \\ \dots & \dots & \dots & \dots & \dots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \phi_2(\mathbf{x}_N) & \dots & \phi_m(\mathbf{x}_N) \end{bmatrix}$$

4



Polynomial regression

- Consider 1D for simplicity:

$$f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_mx^m.$$

- No longer linear in x – but still linear in \mathbf{w} !
- Define $\phi(x) = [1, x, x^2, \dots, x^m]^T$
- Then, $f(x; \mathbf{w}) = \mathbf{w} \cdot \phi(x)$ and we are back to the familiar simple linear regression. The least squares solution:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \text{ where } \mathbf{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \\ 1 & x_2 & x_2^2 & \dots & x_2^m \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_N & x_N^2 & \dots & x_N^m \end{bmatrix}$$

5



General additive regression models

- A general extension of the linear regression model:

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1\phi_1(\mathbf{x}) + w_2\phi_2(\mathbf{x}) + \dots + w_m\phi_m(\mathbf{x}),$$

where $\phi_j(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}, j = 1, \dots, m$ are the *basis functions*.

- This is still linear in \mathbf{w} ,

$$f(\mathbf{x}; \mathbf{w}) = \mathbf{w} \cdot \phi(\mathbf{x})$$

even when ϕ is non-linear in the inputs \mathbf{x} .

6



General additive regression models

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1\phi_1(\mathbf{x}) + w_2\phi_2(\mathbf{x}) + \dots + w_m\phi_m(\mathbf{x}),$$

- Still the same ML estimation technique applies:

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

where \mathbf{X} is the *design matrix*

$$\begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \dots & \phi_m(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \phi_2(\mathbf{x}_2) & \dots & \phi_m(\mathbf{x}_2) \\ \dots & \dots & \dots & \dots & \dots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \phi_2(\mathbf{x}_N) & \dots & \phi_m(\mathbf{x}_N) \end{bmatrix}$$

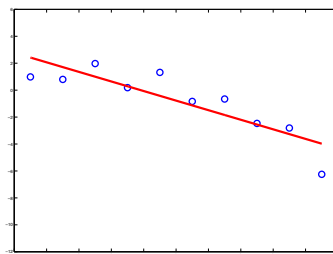
(for convenience we will denote $\phi_0(\mathbf{x}) \equiv 1$)

7

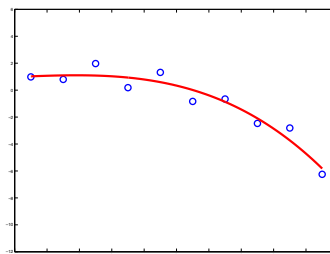


Model complexity and overfitting

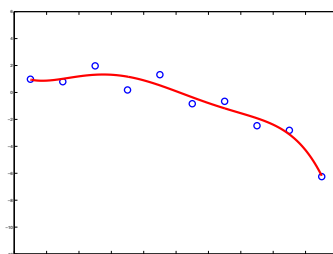
- Data drawn from 3rd order model:



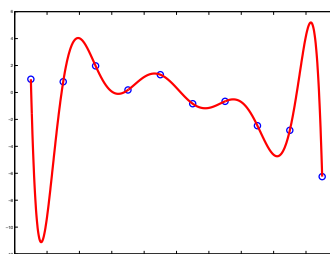
$m = 1$



$m = 3$



$m = 5$



$m = 10$

8



How to avoid overfitting

- The basic idea: if a model overfits (is *too sensitive* to data) it will be unstable. I.e. removal part of the data will change the fit significantly.
- We can *hold out* part of the data; this is *validation (val)* set, or *development (dev)* set.
- Fit the model to the rest, and then test on the heldout data.
- What are the problems of this approach?
 - If the heldout set too small, we are susceptible to chance.
 - If it's too large, we get overly pessimistic (training on too little data compared to what we could do).

9



Cross-validation

- The improved holdout method: *k-fold cross-validation*
 - Partition data into *k* roughly equal parts;
 - Train on all but *j*-th part, test on *j*-th part



10



Cross-validation

- The improved holdout method: *k-fold cross-validation*
 - Partition data into k roughly equal parts;
 - Train on all but j -th part, test on j -th part



11



Cross-validation

- The improved holdout method: *k-fold cross-validation*
 - Partition data into k roughly equal parts;
 - Train on all but j -th part, test on j -th part



12



Cross-validation

- The improved holdout method: *k-fold cross-validation*
 - Partition data into k roughly equal parts;
 - Train on all but j -th part, **test** on j -th part



- An extreme case: *leave-one-out cross-validation*

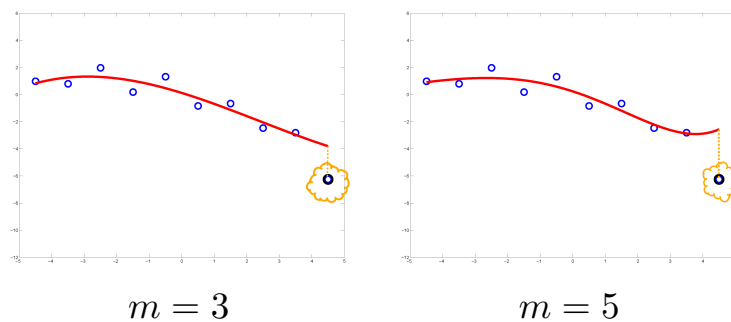
$$\hat{L}_{cv} = \frac{1}{N} \sum_{i=1}^N (y_i - f(\mathbf{x}_i; \hat{\mathbf{w}}_{-i}))^2$$

where $\hat{\mathbf{w}}_{-i}$ is fit to all the data but the i -th example.

13



Cross-validation: example

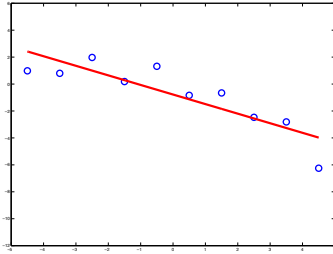


- This is a very good estimate, although expensive to compute
 - Need to run N estimation problems each on $N - 1$ examples!
 - An important research area: devising tricks for efficiently computing cross-validation estimates (by taking advantage of overlap between folds).

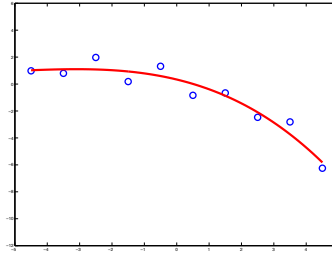
14



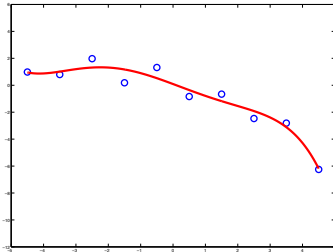
Cross-validation: example



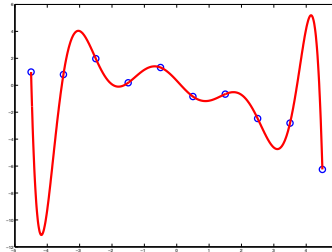
$$m = 1 : L = 1.4, \hat{L}_{cv} = 2.6$$



$$m = 3 : L = 0.4, \hat{L}_{cv} = 1.3$$



$$m = 5 : L = 0.3, \hat{L}_{cv} = 2.7$$



$$m = 10 : L = 0, \hat{L}_{cv} = 4 \times 10^4$$

15



Understanding overfitting

- Cross validation provides some means of dealing with overfitting
- What is the source of overfitting? Why do some models overfit more than others?
- We can try to get some insight by thinking about the estimation process for model parameters

16



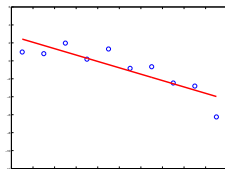
Roadmap

- So far: least squares regression (with arbitrary feature functions)
 - Closed form solution for maximum likelihood
 - Overfitting is a problem
- Today: regularization – main tool to combat overfitting
- Also: gradient descent as an alternative to closed form solution

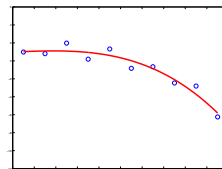
17



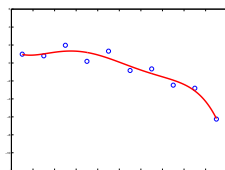
Controlling for overfitting



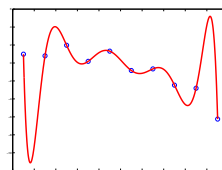
$$L = 1.4, \hat{L}_{cv} = 2.6$$



$$L = 0.4, \hat{L}_{cv} = 1.3$$



$$L = 0.3, \hat{L}_{cv} = 2.7$$



$$L = 0, \hat{L}_{cv} = 4 \times 10^4$$

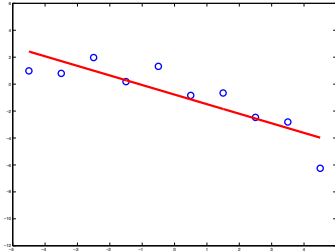
- More complex model (10th degree) overfits more than simple model (linear)
- Pure ERM would always prefer complex models
- Holdout/validation/cross-validation is a way to control for this in *model selection*

18

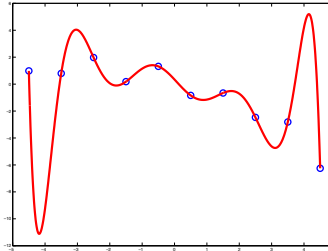


Model complexity - intuition

- Intuitively, the complexity of the model can be measured by the number of “degrees of freedom” (independent parameters).
 - The more complex the model, the more data needed to fit
 \Rightarrow For a given number of points, a more complex model more likely to overfit.



$m = 1, 2$ parameters



$m = 10, 11$ parameters

- This is an issue only because of finite training data!

19



Penalizing model complexity

- Idea 1: restrict model complexity based on amount of data
 - Rule of thumb: ≈ 10 examples per parameter
- Idea 2: directly penalize by the number of parameters.
 Akaike information criterion (AIC): maximize

$$\log p(X | \hat{\mathbf{w}}) - \#\text{params}$$

- But: Definition of model complexity as a number of parameters is a bit too simplistic. Consider feature vector

$$\phi x = \begin{bmatrix} 1 & x & -2x & 2x & x^2 & \frac{1}{2}x^2 \end{bmatrix}$$

Does linear regression $\phi(x) \rightarrow y$ really have 6 parameters?

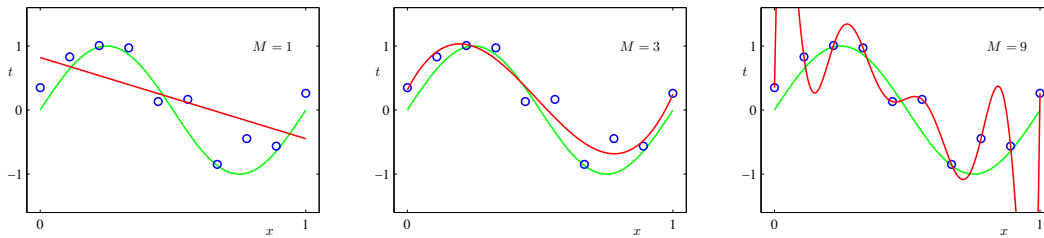
- Idea: look at the behavior of the values of \mathbf{w}^*

20



Linear regression complexity

- Example: polynomial regression, true [from Bishop, Ch. 1]



- Value of the optimal (ML) regression 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

21



Description length

- Intuition: should penalize not the parameters, but the number of bits required to encode the parameters
- With finite set of parameter values, these are equivalent
- With “infinite” set, we can limit the effective number of degrees of freedom by restricting the value of the parameters.
- Then we have penalized log-likelihood:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} \left\{ \frac{1}{2} \sum_{i=1}^N \log p(\text{data}_i; \mathbf{w}) - \text{penalty}(\mathbf{w}) \right\}$$

22



Shrinkage methods

- Shrinkage methods impose penalty on the size of \mathbf{w}
- We can measure “size” in a few different ways. Let us start with L_2 norm:

$$\mathbf{w}_{\text{ridge}}^* = \underset{\mathbf{w}}{\operatorname{argmax}} \left\{ \sum_{i=1}^N \log p(\text{data}_i; \mathbf{w}) - \lambda \|\mathbf{w}\|^2 \right\}$$

in regression “data_{*i*}” = $y_i | \mathbf{x}_i$

- This is *ridge regression*; λ is the *regularization* parameter
- Does it matter that log-likelihood is not averaged? Consider relative effect of the value of λ

23



Ridge regression

$$\mathbf{w}_{\text{ridge}}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 + \lambda \sum_{j=1}^m w_j^2 \right\}$$

- Recall: $\mathbf{w} = [w_0, w_1, \dots, w_m]$
- Usually do not include w_0 in regularization (why?)
- Closed form solution:

$$\hat{\mathbf{w}}_{\text{ridge}}^* = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

- Careful: solution *not* invariant to scaling! Should normalize input before solving.

24



Lasso regression

- The L_1 -penalized maximum likelihood under Gaussian noise model:

$$\mathbf{w}_{\text{lasso}}^* = \underset{\mathbf{w}}{\operatorname{argmax}} \left\{ - \sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 - \lambda \sum_{j=1}^m |w_j| \right\}$$

- This is still concave (i.e. unique maximum), but not “smooth” (differentiable).
- Can solve it efficiently using convex programming methods or first-order numerical optimization (gradient descent)
- Why is it called “lasso”?
least absolute shrinkage and selection operator

25



Optimization of ridge regression

- Can rewrite the optimization problem

$$\min_{\mathbf{w}} \sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 + \lambda \sum_{j=1}^m w_j^2$$

in the proper objective/constraint form:

$$\begin{aligned} & \min_{\mathbf{w}} \sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 \\ & \text{subject to } \sum_{j=1}^m w_j^2 \leq t \end{aligned}$$

- Correspondence $\lambda \Rightarrow t$ can be shown using Lagrange multipliers.

26



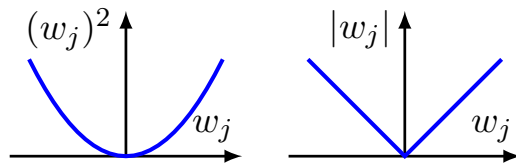
Optimization for Lasso

- Similarly, for Lasso:

$$\min_{\mathbf{w}} \sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2$$

$$\text{subject to } \sum_{j=1}^m |w_j| \leq t$$

- Compare shape of the penalty as a function of w_j :



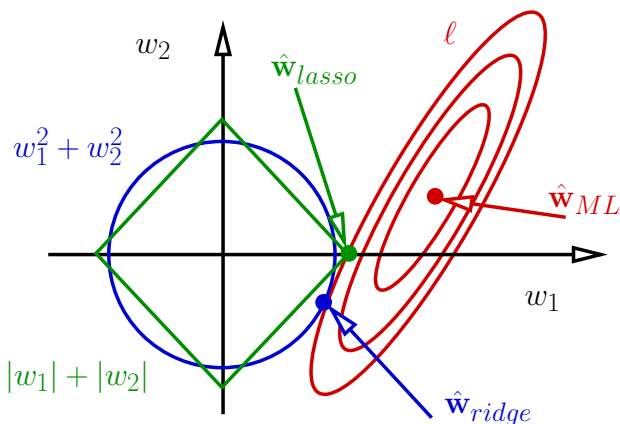
27



Lasso vs. ridge: geometry of error surfaces

- An equivalent formulation for L_p regularization: constrained maximization

$$\hat{\mathbf{w}} = \underset{\mathbf{w}: \sum_{j=1}^m |w_j|^p \leq \beta}{\operatorname{argmax}} - \sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2.$$



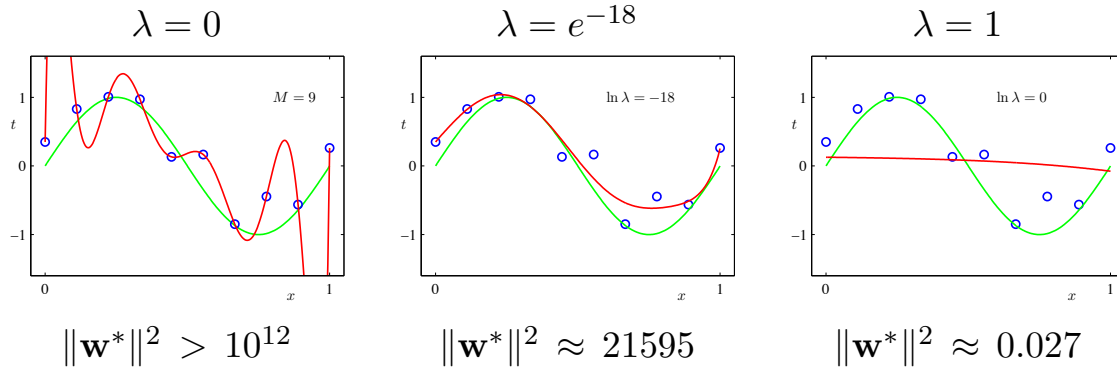
- With sufficiently large λ (=sufficiently small β) lasso leads to *sparsity*.
- Must explicitly solve the above optimization problem – e.g., using Lagrange multipliers.

28



Choice of λ

- Example [from Bishop, Ch. 1]: 9th deg polynomial with varying λ :



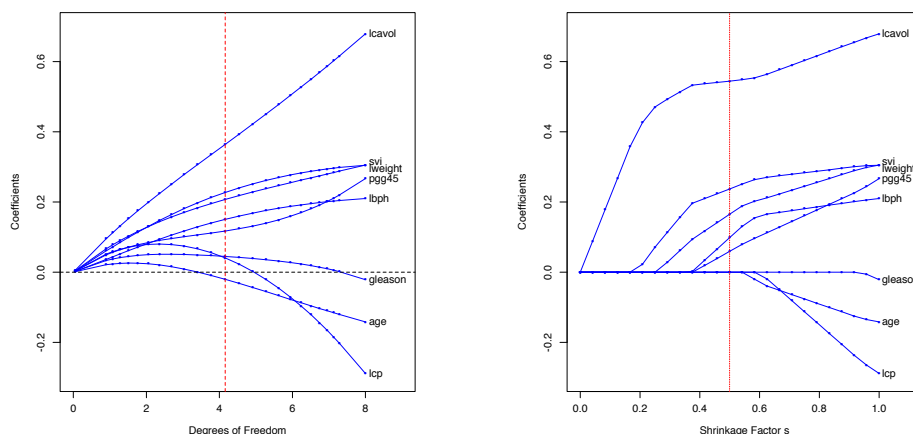
- Most often: choose λ by (cross) validation

29



Example: lasso vs. ridge regularization paths

- Example: prostate data [Hastie, Tibshirani and Friedman]
Red lines: choice of λ by 10-fold CV.

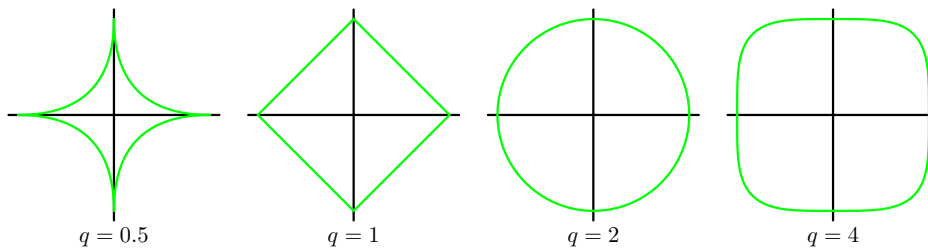


30



General view of L_q penalty

- Can be creative in design of penalty function $\|\mathbf{w}\|_q$



- For $q > 1$, no sparsity is achieved.
- For $q < 1$, non-convex
- What about L_0 ?

$$\min_{\mathbf{w}} \sum (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 \quad \text{s.t.} \quad |\{w_j : w_j > 0\}| \leq M$$

is NP-hard

31

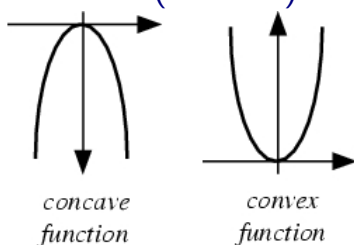


Beyond closed form solution

- So far: solve (least squares) regression with a closed form solution

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Sometimes we can not do this. E.g., the data matrix is too large to compute the pseudoinverse for
- If we move away from simple squared loss (e.g., in PS1: asymmetric loss) also lose the closed form solution
- Alternative: numerical optimization – gradient descent
- Consider (for now) convex or concave functions

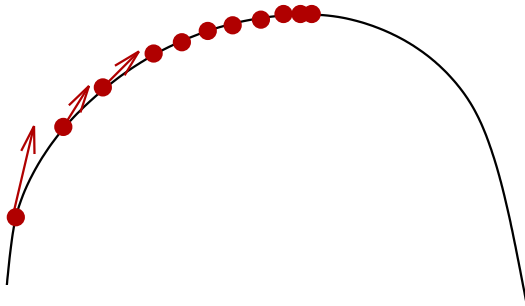


32



Gradient ascent/descent

- The idea behind gradient ascent: “hill climbing” on the function surface.



- Start at a (random) location
- Make steps in the direction of maximal altitude increase.

- An equivalent: gradient *descent* on the *convex* loss $-\log p(y | \mathbf{x}; \mathbf{w})$

33



Gradient descent algorithm on $f(\mathbf{X}, \mathbf{y}; \mathbf{w})$

- Iteration counter $t = 0$
- Initialize $\mathbf{w}^{(t)}$ (to zero or a small random vector)
- for $t = 1, \dots$:
 compute gradient

$$\mathbf{g}^{(t)} = \nabla f(\mathbf{X}, \mathbf{y}; \mathbf{w}^{(t-1)})$$

update model

$$\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} - \eta \mathbf{g}^{(t)}$$

check for convergence

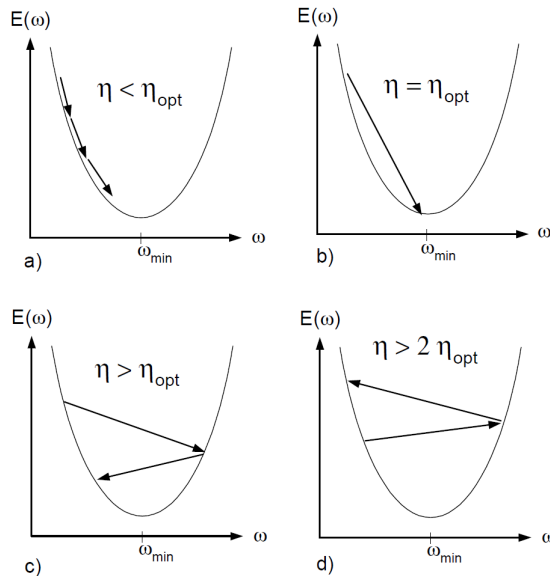
- The *learning rate* η controls the step size

34



Gradient descent convergence

- Generally, for convex functions, gradient descent will converge
- Setting the learning rate η may be very important to ensure rapid convergence



From Lecun et al, 1996

35



Gradient descent convergence

- A lot of theory on convergence of gradient descent
- Usually relies on various properties of the objective function: strong convexity, smoothness, etc.
- In practice, need to monitor the objective, tweak learning rate, and consider stopping (“convergence”) criteria
- Common criteria (often use a combination):
 - Maximum number of iterations (time budget)
 - Minimum required change in objective value (loss)
 - Minimum required change in model parameters (\mathbf{w})
- If stopped because of max iterations: may not have converged
- Problematic criteria: monitor absolute (not relative) value of something like objective or parameters. Often hard to know what the “right” value for these is.

36

