

EdX and its Members use cookies and other tracking technologies for performance, analytics, and marketing purposes. By using this website, you accept this use. Learn more about these technologies in the [Privacy Policy](#).



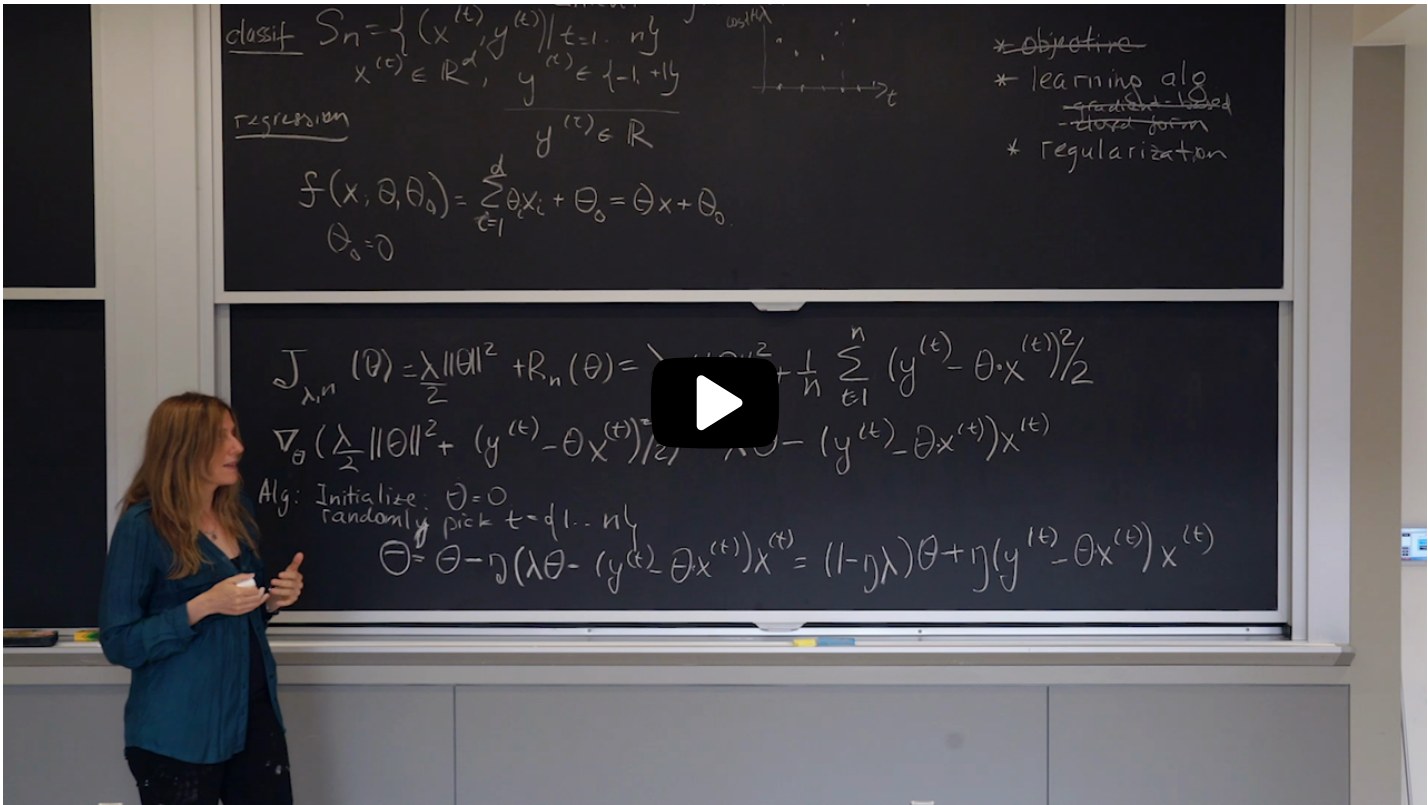
[Unit 2 Nonlinear Classification](#),  
[Linear regression, Collaborative](#)

[Course](#) > [Filtering \(2 weeks\)](#) > [Lecture 5. Linear Regression](#) > 9. Closing Comment

## 9. Closing Comment

### Closing Comment

[Start of transcript. Skip to the end.](#)



Now what I want to do before we close today's lecture is actually is to say jointly what this regularization is doing. It doesn't matter how, at this point, which algorithm do you use. I want to bring you back to this formula, to the Suivche regression formula and think together with me, what does it do? Like previously, when we had our normal



**Video**  
[Download video file](#)

**Transcripts**  
[Download SubRip \(.srt\) file](#)

[Download Text \(.txt\) file](#)

(Optional) Equivalence of regularization to a Gaussian Prior on Weights

Show

Discussion

Hide Discussion

**Topic:** Unit 2 Nonlinear Classification, Linear regression, Collaborative Filtering (2 weeks):Lecture 5. Linear Regression / 9. Closing Comment

Add a Post

◀ All Posts

Derivation of loglikelihood inside, spoiler alert

discussion posted a day ago by [Cool7](#) (Community TA)

As title. This is easier than last one. Just put it here in case somebody interested. I'm practicing my latex writing, lol.

$$\begin{aligned} &\log\left(\prod_{t=1}^n \mathcal{N}(y_t|\theta x_t, \sigma^2) \mathcal{N}(\theta|0, \lambda^{-1})\right) \\ &= \sum_{t=1}^n (\log(\mathcal{N}(y_t|\theta x_t, \sigma^2)) + \log(\mathcal{N}(\theta|0, \lambda^{-1}))) \\ &= n \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + \sum_{t=1}^n \log\left(e^{-\frac{(y_t-\theta x_t)^2}{2\sigma^2}}\right) + n \log\left(\sqrt{\frac{\lambda}{2\pi}}\right) + \sum_{t=1}^n \log\left(e^{-\frac{\lambda\|\theta\|^2}{2}}\right) \\ &= \sum_{t=1}^n \left(-\frac{1}{2\sigma^2}(y_t - \theta x_t)^2 - \frac{\lambda}{2}\|\theta\|^2\right) + n \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) + n \log\left(\sqrt{\frac{\lambda}{2\pi}}\right) \\ &= \sum_{t=1}^n -\frac{1}{2\sigma^2}(y_t - \theta x_t)^2 - \frac{1}{2}\lambda\|\theta\|^2 + \text{constant} \end{aligned}$$

My understanding is

- First term is related to posterior distribution, it represents the accuracy of the estimation/training loss/bias.
- Second term is related to prior distribution, it represents the regularization(recall we imposed it on) / variance.

Thus  $\lambda$  as hyper parameter is to adjust the weights between bias and variance, inline with the error decomposition discussed a few pages before.

This post is visible to everyone.



Add a Response

**Alexander Konstantinidis**  
about 5 hours ago



Another way to view this, is to consider  $\lambda$  as expressing the degree of our certainty (prior belief) that there is no real explanatory value in the model or stated differently very few if any of the predictors truly matter. (This is because  $\lambda$  is the inverse of the variance of probabilistic theta). The higher the  $\lambda$  the more evidence will be required to arrive to a complex model and vice versa.

Indeed, this is a very interesting interpretation. In the extreme case, where lambda is infinity, it means your prior belief is so strong that no matter what data is presented, the hard coded parameters do not change. On the other extreme, when lambda is 0, variance is infinity and thus you don't have a prior belief. Data takes control of everything, even if there're a lot of noise. So a moderate lambda lets the model to learn from data, but regularizes the parameters so that do not deviate too much from the prior belief.

posted about 5 hours ago by **FutureStar**

Add a comment

Showing all responses

Add a response:

Preview

Submit

Learn About Verified Certificates