



## Naive Bayes

Naive Bayes Classification یک گروه از دسته‌بندی‌های بر پایه احتمالات است که به صورت زیر عمل می‌کند.

$$\operatorname{argmax}_k \left\{ \underbrace{P(y = k|x)}_{\text{MAP Rule}} \right\}$$

$$P(y = k|x) = \frac{P(x|y = k) P(y = k)}{P(x)}$$

در این شرایط تخمین عبارت  $P(y = k)$  ساده است اما تخمین عبارت  $P(x|y = k)$  به دلیل ابعاد بالای ماتریس سخت است. برای ساده‌سازی این عبارت اگر مقادیر  $x$  از هم مستقل باشند می‌توان عبارت بالا را به صورت زیر نوشت:

$$P(x|y = k) = \prod P(x_i|y = k)$$

لازم به ذکر است که اگر  $y$  جزو داده‌ها باشد،  $x$  ها از یکدیگر مستقل هستند. در غیر این صورت مستقل نخواهند بود و وابسته‌اند چون در یک سمت مشترکند و وابستگی دارند. در نهایت گروه‌بندی بر اساس  $y$  های متفاوت انجام می‌گیرد.

$$\hat{P}(x_1|y = 1), \hat{P}(x_1|y = 2), \dots, \hat{P}(x_1|y = n)$$

## Logistic Regression

در این قسمت می‌خواهیم جواب نهایی ( $P(y = k|x)$ ) را حدس بزنیم.

### دو کلاسه:

فرض می‌کنیم که مرز بین نواحی خطی است و با یک hyperplane می‌توان فضا را به دو قسمت تبدیل کرد. می‌خواهیم  $P(y = 1|x) \leq P(y = 2|x)$  مقایسه کنیم.

$$\frac{P(y = 1|x)}{P(y = 2|x)} \leq 1 \implies \log \left( \frac{P(y = 1|x)}{P(y = 2|x)} \right) \leq \log(1)$$

اگر حاصل لگاریتم اول برابر با  $\beta^T x$  بود، می‌توانستیم این مقدار را با صفر مقایسه کنیم و در نتیجه عبارت بالا را به صورت زیر بازنویسی می‌کنیم.

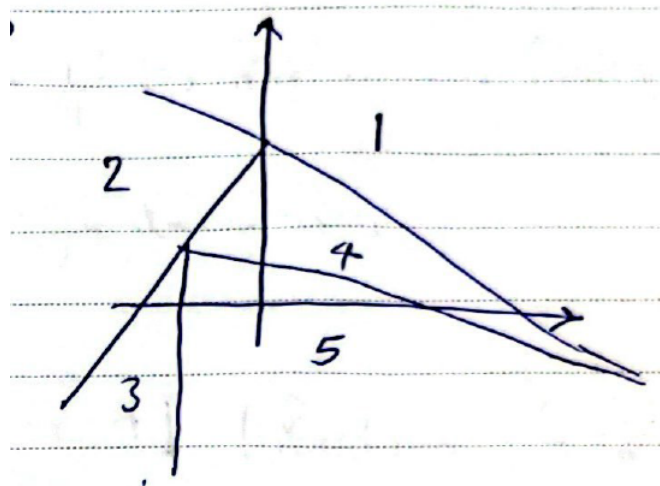
$$P(y = 1|x) = \frac{e^{\beta^T x}}{1 + e^{\beta^T x}}, \quad P(y = 2|x) = \frac{1}{1 + e^{\beta^T x}}$$

حال با استفاده از تخمینگر MLE مقدار  $\beta$  را تخمین می‌زنیم.

$$\begin{aligned} P(D|\beta) &= \prod P(y_i|x_i, \beta) \implies \sum \log(P(y_i|x_i, \beta)) = \\ &= \sum \log \left( \frac{e^{\beta^T x^{(i)}}}{1 + e^{\beta^T x^{(i)}}} \right) + \sum \log \left( \frac{1}{1 + e^{\beta^T x^{(i)}}} \right) \\ &\implies \sum \left\{ y_i \log \left( \frac{e^{\beta^T x^{(i)}}}{1 + e^{\beta^T x^{(i)}}} \right) + (1 - y_i) \log \left( \frac{1}{1 + e^{\beta^T x^{(i)}}} \right) \right\} \end{aligned}$$

با استفاده از مشتق، می‌توان به جواب بهینه برای مسئله بالا رسید.

## چندکلاسه



برای یک مسئله  $k$  کلاسه نیاز به  $k - 1$  hyperplane داریم که بتوان فضا را به  $k$  قسمت تقسیم کرد. با مقایسه کردن هر دو فضا با یکدیگر می‌خواهیم ضرایب  $\beta$  را بدست آوریم.

$$P(y = i|x) \leq P(y = j|x) \implies \log \left( \frac{P(y = i|x)}{P(y = j|x)} \right) \leq 0 \implies \beta_{ij}^T x \leq 0$$

اگر در این مرحله  $y = i$  را با  $y = k$  مقایسه کنیم، این مدل consistent می‌باشد. یا به عبارتی دیگر یک گروه را مرجع در نظر بگیریم و بقیه را نسبت به آن مقایسه کنیم.

$$P(y = 1|x) = \frac{e^{\beta_1 x}}{z}, P(y = 2|x) = \frac{e^{\beta_2 x}}{z}, \dots, P(y = n|x) = \frac{1}{z}$$

$$\text{Where } z = 1 + \sum_{i=1}^{k-1} e^{\beta_i x}$$

سوال: اگر ۱۰ گروه و ۱۰۰ بعد داشته باشیم. چند پارامتر خواهیم داشت؟  
۹ hyperplane خواهیم داشت که هر کدام به ۱۰۱ پارامتر نیاز خواهند داشت.

$$9 \times 101$$

تعریف ۱ (آنتروپی).

$$H((p_1, \dots, p_k)) = - \sum p_i \log(p_i)$$

تعریف ۲ (دیورژانس).

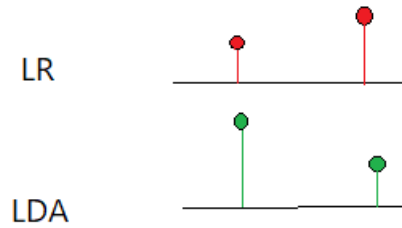
$$D((p_1, \dots, p_k) || (q_1, \dots, q_k)) = \sum p_i \log\left(\frac{p_i}{q_i}\right)$$

## مقایسه مدل‌ها

ارزیابی مدل یک مرحله مهم در تعیین عملکرد بهینه و همچنین یک معیار تاثیرگذار در انتخاب مدل مورد نظر می‌باشد. به طور مثال فرض کنید که برای یک ست داده می‌خواهیم تصمیم بگیریم که از کدام یک از دو روش LDA و یا Logistic Regression استفاده کنیم. اگر فرض کنیم که حجم داده‌ها به اندازه کافی زیاد است، می‌توان به صورت زیر عمل کرد.

ابتدا داده‌ها را به ۴ قسمت تقسیم می‌کنیم. سپس از  $\frac{1}{4}$  داده‌ها به ترتیب برای *train* کردن مدل‌های LDA و Logistic Regression (توجه کنید که برای هر مدل از یک دسته داده مجزا استفاده می‌شود) و از  $\frac{1}{4}$  داده‌ها برای *test* کردن این مدل‌ها استفاده می‌کنیم. (در نظر داشته باشید که بین این ۴ دسته از داده‌ها اشتراکی وجود ندارد. لزوماً مدلی که در *train* نتیجه بهتری داشته باشد مدل بهتری نیست.)

آنها را با یکدیگر مقایسه کنیم. برای مثال اگر نتیجه حاصل مشابه زیر باشد ما مدل LDA را انتخاب می‌کنیم زیرا می‌توان مشاهده کرد که مدل دوم در واقع دارای پیچیدگی بیشتری نسبت به پیچیدگی داده‌های ما هست و عملاً خطای کمتر *train* گمراه کننده است.



شکل ۱: مثال نتیجه مقایسه دو مدل

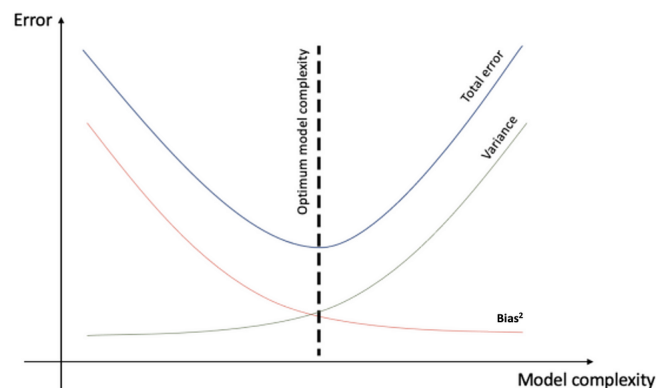
می‌توان از تعمیم دادن این ایده به تعداد بیشتری مدل استفاده کرد. در حالت کلی با توجه به ضرایب  $\beta$  و تعداد پارامترهای به کار رفته می‌توانیم بی‌نهایت مدل مختلف داشته باشیم و خطای مدل‌ها را به صورت زیر محاسبه کنیم:

$$M_1 \rightarrow \frac{1}{n} \sum L(y_i, f_1(x_i))$$

$$M_2 \rightarrow \frac{1}{n} \sum L(y_i, f_2(x_i))$$

⋮

این روش زمانی به پاسخ صحیح خواهد رسید و برای تمامی  $M$  ها تخمین خوبی ارائه خواهد داد که تعداد داده‌ها زیاد باشد. نموداری از خطای  $train$  و  $test$  بر حسب  $complexity$  تشکیل می‌دهیم. بنابراین در نهایت مدلی مد نظر ما است که هم در  $train$  و هم در  $test$  خطای کمی داشته باشد و عملکرد بهتر در تنها یکی از این بخش‌ها کافی نیست. در این مقایسه ممکن است پیچیدگی مدل در مقابل عملکرد آن قرار بگیرد. با افزایش پیچیدگی مدل می‌توان در داده‌های  $train$  به نتایج بهتری دست یافت (خطای بایاس را کاهش می‌دهیم) اما همین افزایش پیچیدگی می‌تواند در داده‌های  $test$  به نتایج ضعیف‌تر منجر شود (خطای واریانس افزایش می‌یابد). در حقیقت در چنین حالتی، ما از مدلی استفاده کرده‌ایم که از مدل واقعی پیچیده‌تر بوده و به اصطلاح  $Overfit$  رخ داده است.



شکل ۲: نمودار دقت مدل با توجه به پیچیدگی آن

به همین منظور واجب است که در مدل خود پیچیدگی را کنترل کنیم. به متغیری که از آن به منظور کنترل پیچیدگی مدل استفاده می‌کنیم  $hyperparameter$  می‌گویند و آن را با نماد  $\alpha$  نمایش می‌دهیم. به طور مثال مدل

خطی زیر را در نظر بگیرید:

$$f_{\beta} = \sum_{i=1}^p \beta_i X_i$$

$$s.t. \quad \sum_{i=1}^p \beta_i^2 \leq \alpha$$

در اینجا با استفاده از متغیر  $\alpha$  می‌توانیم پیچیدگی مدل را کنترل کنیم به طوری که با افزایش  $\alpha$  پیچیدگی افزایش و با کاهش آن، کاهش می‌یابد. توجه کنید که برای داده‌های *typical* مقدار  $\alpha$  تک پارامتری است و می‌توان از تعدادی از نمونه‌ها به عنوان داده *train* استفاده کرد، در حالی که اگر ابعاد  $\alpha$  زیاد باشد، ممکن است چنین کاری عملی نباشد.

## Variance and Bias

$$x : F(x)$$

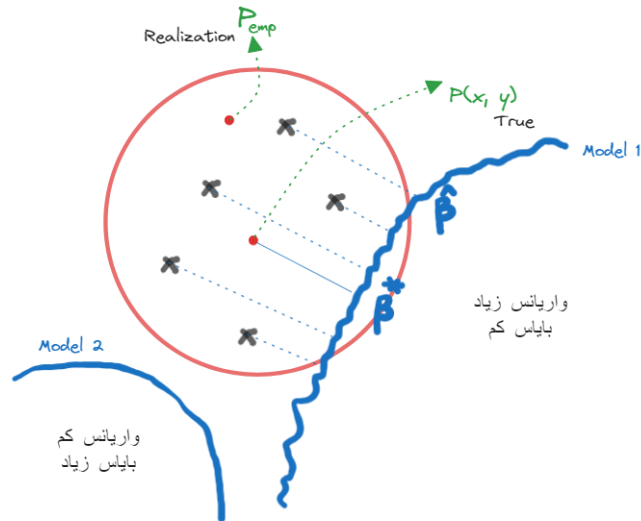
$$F_{emp}(x) = \frac{\sum_{i=1}^m \mathbb{1}[x_i \leq x]}{m}$$

$$f_{emp}(x) = \frac{1}{m} \sum_{i=1}^m \delta(x - x_i)$$

حال سوال مهمی پیش می‌آید، اینکه اگر قرار باشد resampling صورت بگیرد و دوباره همانقدر قبلی برداریم چند درصد داده‌های جدید با قبلی یکسان خواهد بود؟  
مشاهده میشود که با میل کردن *sample size* به بی‌نهایت، مقدار عبارت زیر به  $0.63m$  میل میکند که نشان می‌دهد بیش از نصف داده‌ها تکراری خواهند بود.

$$1 - \left(\frac{m-1}{m}\right)^m \rightarrow \lim_{m \rightarrow \infty} \left(1 - \left(\frac{m-1}{m}\right)^m\right) = m(1 - e^{-1}) \approx 0.63m$$

$$P_{emp} : D = \{(x_i, y_i)\}$$



$$Predict : \hat{Y} = f_{\beta}(x)$$

هدف:

$$\min_{\hat{Y}=f_{\beta}(x)} E(L(Y, \hat{Y}))$$

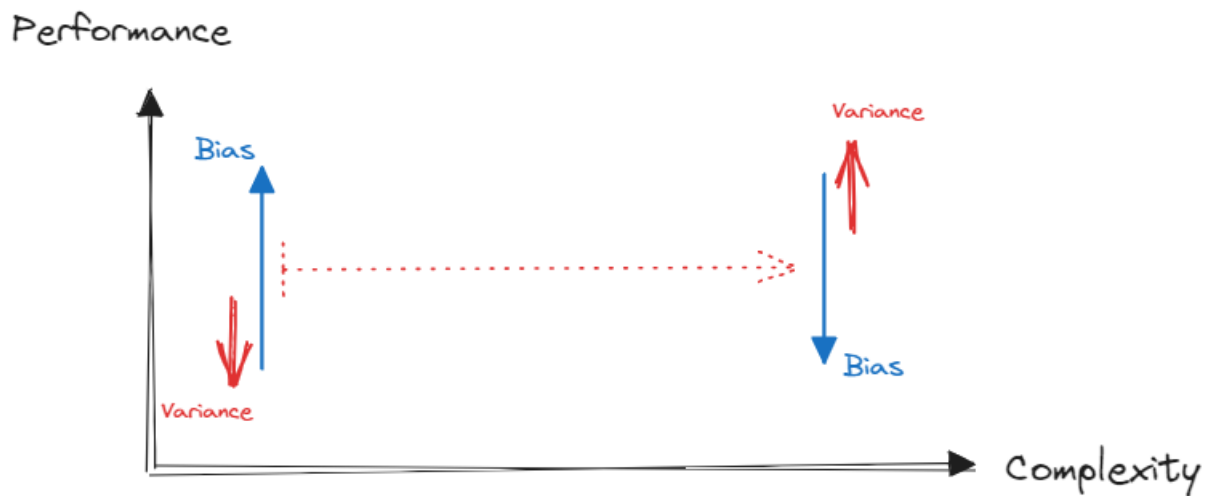
حالا یک  $Err_T$  هم داریم که نشان دهنده میزان خطا بر روی داده های آموزش می باشد.

$$\begin{aligned} Err_T &= E[(Y - \hat{Y})^2 | T] = E[(f(x) + \epsilon - f_{\beta}(x))^2 | T] \\ &= E[\epsilon + f(x) - \bar{f}_{\beta}(x) + \bar{f}_{\beta}(x) - f_{\beta}(x)]^2 | T \\ &= \delta_E^2 + bias^2 + Var \end{aligned}$$

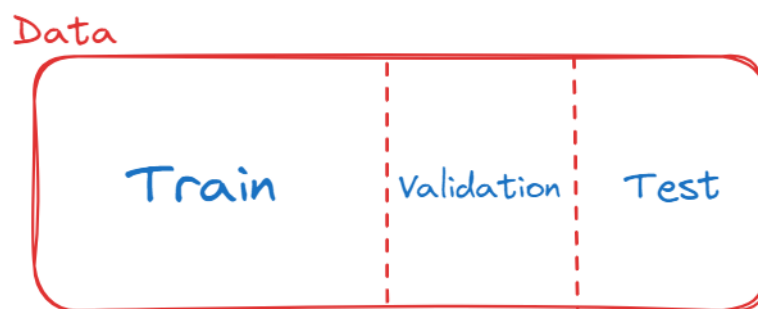
بایاس نشان دهنده فاصله ایست که تا خروجی اصلی داریم.  
و واریانس مقدار آزادی و حرکتی که مدل ما می تواند انجام بدهد.

## Model Complexity

اگر پیچیدگی مدل ما زیاد باشد، ما واریانس زیاد و بایاس کمی خواهیم داشت. اگر ما سعی کنیم تا مدل را **simplify** کنیم یعنی از پیچیدگی آن بکاهیم، به نظر می آید که واریانس را کمتر و بایاس را زیادتر کرده ایم؛ البته باید توجه داشته باشیم که همیشه اینگونه نیست.  
پس در نتیجه معمولا:

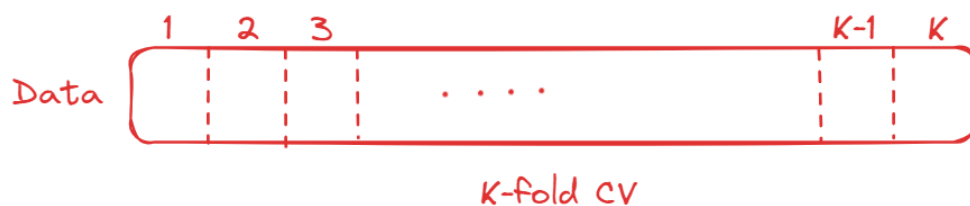


یکی از راه حل هایی که پیش رو داریم تقسیم کردن همه دیتا به قسمت های مختلفی است. بهتر است حدود ۵۰ درصد دیتا را به داده های train، و مابقی آن ها را بین داده های test و validation نصف کنیم. (البته در اهداف و داده های مختلف می توانیم به روش های دیگر این تقسیم را انجام بدهیم)



## Cross Validation

راه قبلی برای وقتی که تعداد داده های زیاد و کافی داشته باشیم مشکلی ندارد و خوب است ولی اگر داده های کمی داشته باشیم باید سراغ راه حل های دیگری برویم.



در این حالت K بار باید عمل training انجام شود و به این شکل، K مدل خواهیم داشت. هر بار یک قسمت را برای validation و بقیه قسمت ها را برای train انتخاب می کنیم.

برای انتخاب مدل نهایی راه های مختلفی وجود دارد که یکی از آن ها را در اینجا می بینیم:

$$f(x) = \frac{1}{K} \sum_{i=1}^K f_{\beta_i}(x)$$

و validation این مدل برابر با:

$$f(x) = \frac{1}{m} \sum_{i=1}^m L(Y_i, \hat{Y}_i)$$

مقدار جواب پیشبینی شده همان نتیجه  $f_{\beta}$  هستش که از روی یک قسمت از داده ها بدست آمده که در آن دفعه برای validation استفاده شده است.

اگر نمونه ای که در هر دفعه برای validation قرار می دهیم فقط یکی باشد به این حالت LOOCV می گوئیم.

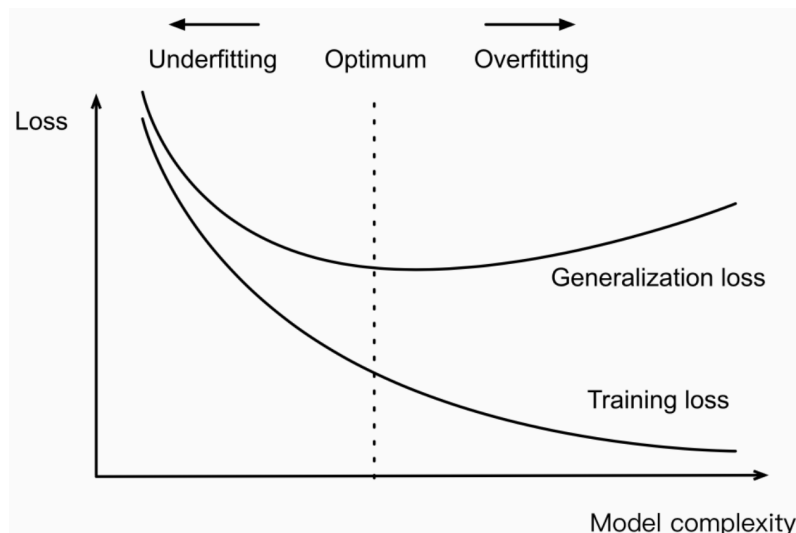
## Regularization

فرض کنید داریم:  $D = \{(x_i, y_i) | i \in \{1, 2, \dots, n\}\}$

مدل ما (همان  $f_{\beta}(x)$ ) می تواند parametric یا non parametric باشد.

همچنین تابع هزینه  $^1$  را اگر با نمادگذاری  $L(y, \hat{y})$  نمایش دهیم، هدف ما پیدا کردن  $\beta$  است که  $\min \sum_{i=1}^n L(y_i, \hat{y}_i)$

را می دهد. منظور از «تعمیم» این است که مدل در داده های جدید که در مجموعه داده آموزشی نبوده (همان داده های تست) عملکرد خوبی داشته باشد. می خواهیم تاثیر complexity مدل را بر روی خطا بر روی داده های تست بررسی کنیم.



همانطور که می بینیم با افزایش complexity خطای prediction نیز افزایش می یابد در نتیجه باید میزان com-plexity مدل را کنترل کنیم و این کار را می توانیم با استفاده از regularization انجام دهیم. Regularization

<sup>1</sup> loss function



می‌تواند explicit یا implicit باشد. در explicit یک ناحیه را در نظر می‌گیریم؛ مانند  $0 \leq \beta_1 \leq 5$  یا  $|\beta_1| < 0.1$  و یا  $\sum |\beta_i| < 0.1$

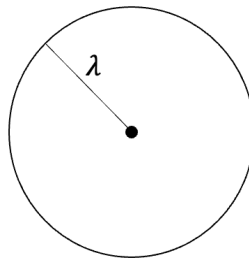
Regularization به صورت implicit در فرایند بهینه‌سازی اعمال می‌شود. در واقع الگوریتم را طوری در نظر می‌گیریم که از یک جا شروع کند و فقط ناحیه اطراف را جستجو کند. در نتیجه همه‌ی فضا در نظر گرفته نشده است. و complexity کمتر شده است. از مثال‌های آن می‌توان به drop out و early stopping اشاره کرد. اگر مدل را بصورت زیر در نظر بگیریم:

$$f_{\beta}(x) = \sum_{i=1}^n \beta_i x_i$$

## انواع Regularization

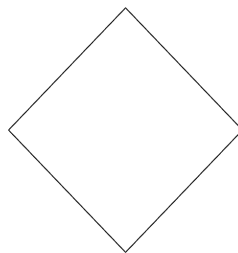
• Ridge Regularization

$$\|\beta\|_2^2 = \sum_{i=1}^n \beta_i^2 \leq \lambda$$



• Lasso Regularization

$$\|\beta\|_1 = \sum_{i=1}^n |\beta_i| \leq \lambda$$



• Subset Selection

$$\|\beta\|_0 = |\{i : \beta_i \neq 0\}| \leq \lambda$$

• Elastic Net

$$\|\beta\|_2^2 \leq \lambda_1$$

$$\|\beta\|_1 \leq \lambda_2$$

Lasso به صورت خودکار یک سری از ضرایب را صفر می‌کند.  
اگر داشته باشیم:

$$\min g(\beta) \\ \text{such that } h(\beta) \leq \lambda$$

می‌توانیم تابع هدف را به صورت زیر بازنویسی کنیم:

$$\min_{\beta} g(\beta) + \mu h(\beta)$$

در این صورت به جای تغییر  $\lambda$  می‌توانیم  $\mu$  ها را تغییر دهیم. چون  $\mu$  ،  $\lambda$  متناظر هم هستند و یعنی به ازای هر  $\mu$  یک  $\lambda$  وجود دارد و برعکس.  
به طور مثال در lasso داریم:

$$\min_{\beta} \frac{1}{n} \sum \|y_i - \beta_i x_i\|^2 \\ \text{s.t } \|\beta\|_1 \leq \lambda$$

که معادلا برابر است با:

$$\min_{\beta} \frac{1}{n} \sum \|y_i - \sum \beta_i x_i\|^2 + \mu \|\beta\|_1$$

هر چه  $\lambda$  بیشتر شود پیچیدگی بیشتر می‌شود اما هر چه  $\mu$  بیشتر شود پیچیدگی کمتر می‌شود.  
در فضای non parametric نیز می‌توان constraint گذاشت:  
به طور مثال:  
بین همه  $f$  ها جستجو شود به طوری که:

$$\text{constrain} = \int f(x)^2 dx \leq \lambda$$

معادلا می‌توان نوشت:

$$\min \sum L(y_i, \hat{y}_i) + \mu \int f(x)^2 dx$$

در شبکه‌های عصبی معمولا از روش‌هایی مانند Gradient Descent ،Drop Out ،Early Stopping استفاده می‌شود و این باعث وجود نوعی Implicit Regularization در آن‌ها می‌گردد.