

# فرآیند Regularization

دوره پایتون و یادگیری ماشین

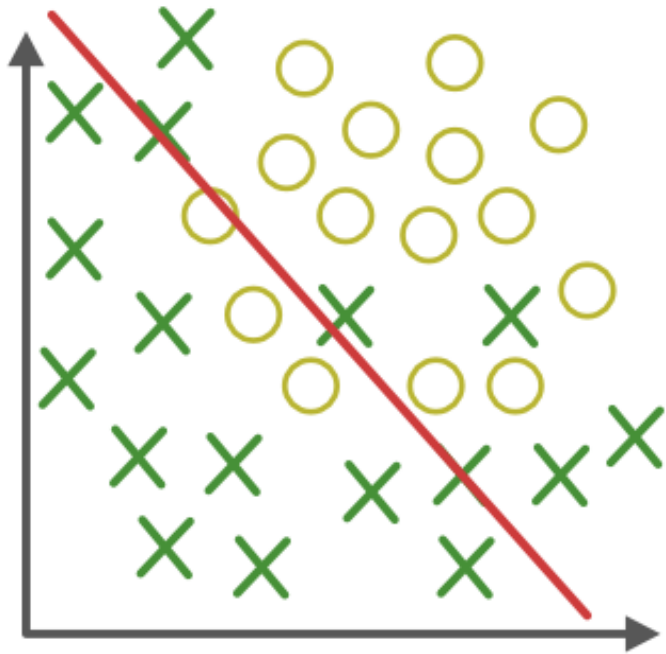


</Monolearn>

# مشکلات مدل یادگیری ماشین

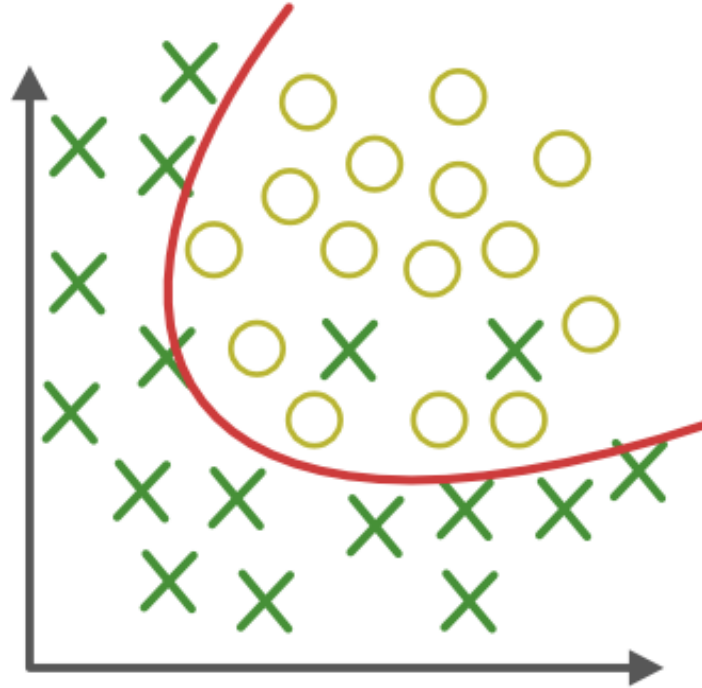
- Overfitting یا بیش برازش: زمانی اتفاق می افتد که مدل یادگیری ماشین، بر روی داده ها به خوبی آموزش ببیند (خطا بر روی داده های train کم و میزان دقت بالا) ولی بر روی داده های تست به خوبی عمل نکند (خطا بر روی داده های test زیاد و میزان دقت پایین).
- Underfitting یا کم برازش: مدل بر روی داده های train به خوبی آموزش نبیند در نتیجه میزان خطا هم بر روی داده ها train و هم test مقدار بالایی خواهد بود و میزان دقت بر روی هر دو مجموعه داده، مقدار کمی خواهد بود. (مدل نه خوب آموزش می بیند و نه خوب عمل یا پیش بینی می کند).

# مشکلات مدل یادگیری ماشین

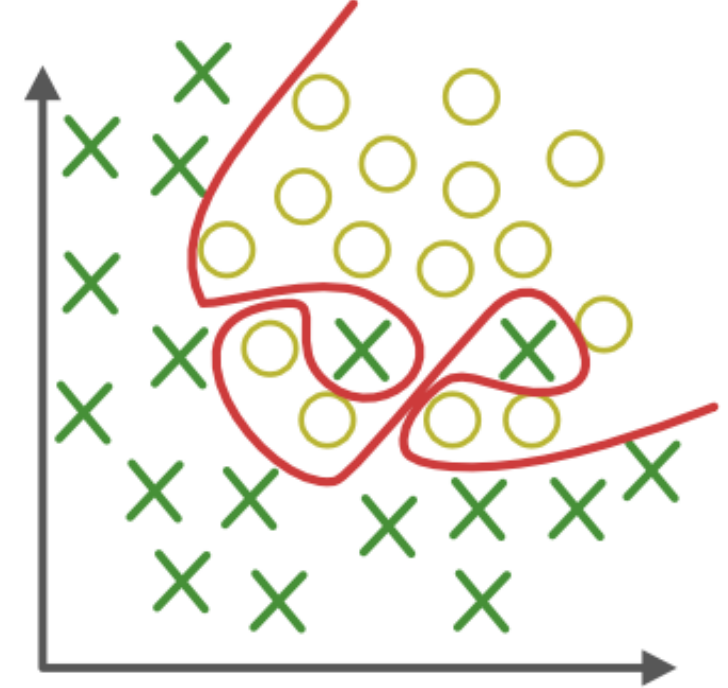


**Under-fitting**

(too simple to  
explain the variance)



**Appropriate-fitting**

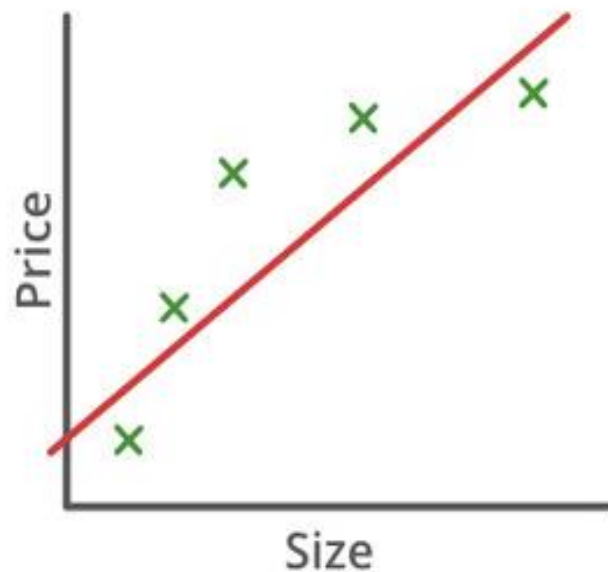


**Over-fitting**

(forcefitting--too  
good to be true)

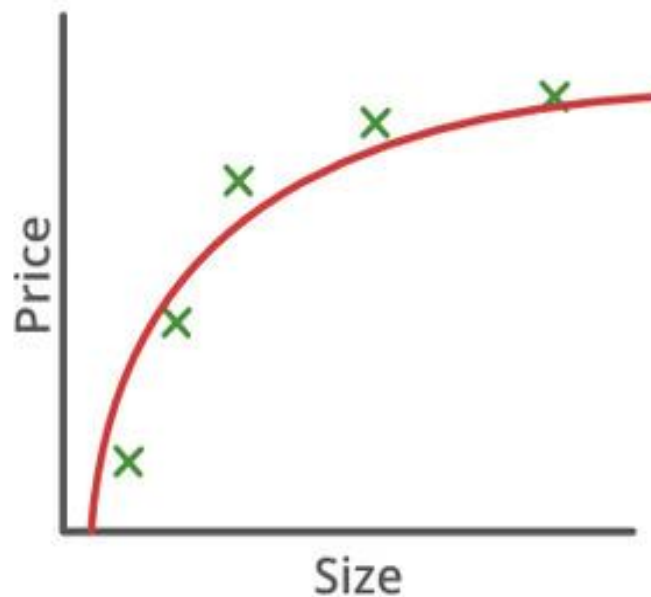


# مشکلات مدل یادگیری ماشین



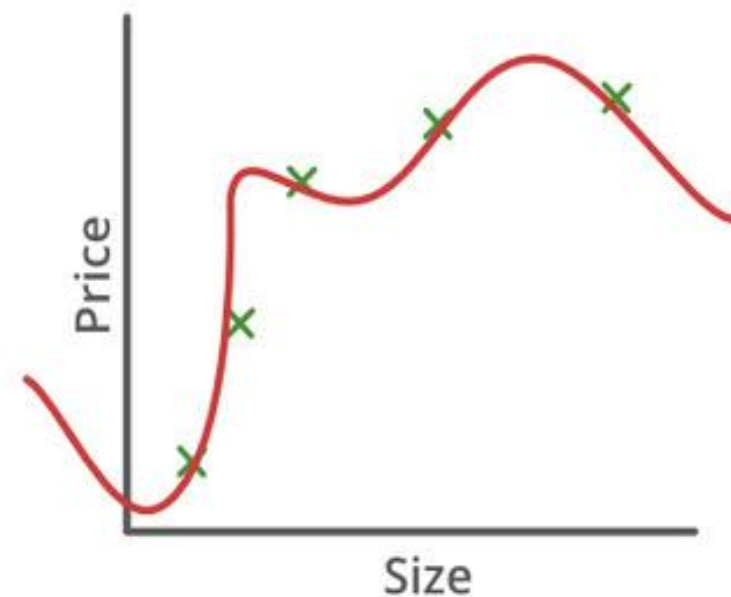
$$\theta_0 + \theta_1 x$$

High bias (underfit)



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

High bias (underfit)



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

High variance  
(overfit)

# Overfitting یا بیش برازش

- زمانی اتفاق می افتد که مدل یادگیری ماشین بر روی داده های آموزشی بیش از حد آموزش ببیند و در نتیجه در تست با مشکل مواجه شود و نتواند خوب پیش بینی کند.
- در این مورد، مدل یادگیری ماشین، نویز موجود در داده ها را نیز یاد می گیرد.

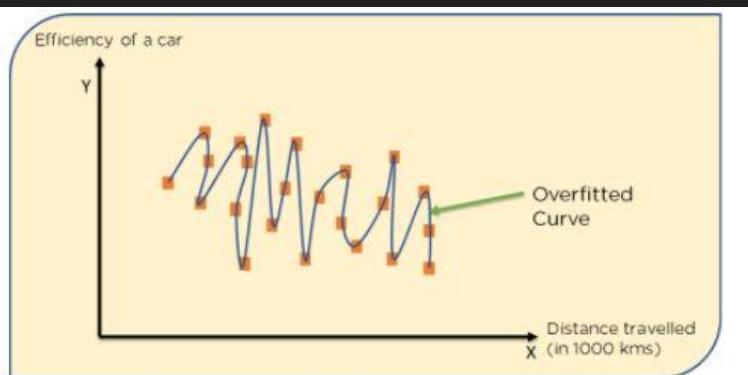
## ○ چه مواقعی اتفاق می افتد؟

- مدل بیش از حد پیچیده است و ویژگی های همخط (Collinear) را در بر می گیرد که واریانس داده ها را افزایش می دهد.

- تعداد ویژگی های داده ها بیشتر یا برابر با تعداد داده هاست (تعداد ویژگی ها بالاست)

- حجم داده ها کم یا بسیار کم است

- داده پیش پردازش نشده، تمیز نیست، نویز دارد

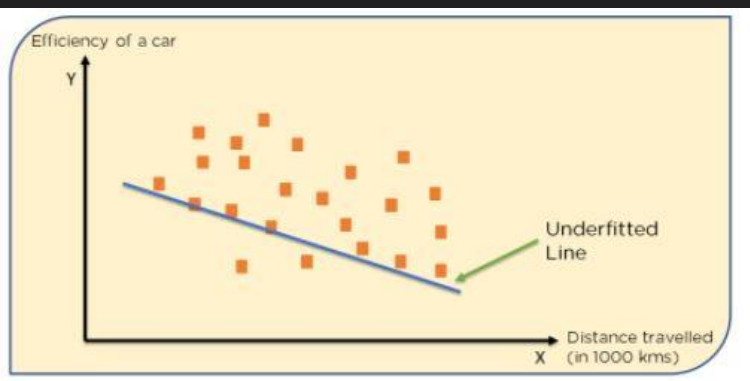


# روش های مقابله با Overfitting

- اعمال روش K-fold cross validation
- افزایش داده ها (جمع آوری داده های بیشتر و همچنین، ایجاد داده های ساختگی)
- انتخاب ویژگی ها به طور موثر
- نظم دهی یا Regularization ( $L_1$  و  $L_2$ ) و (همچنین استفاده از تکنیک Drop Out)
- حذف لایه ها از مدل (از شبکه های عصبی در دیپ لرنینگ)
- توقف زود هنگام عمل یادگیری (در شبکه های عصبی در دیپ لرنینگ)

# Underfitting یا کم برازش

- زمانی اتفاق می افتد که مدل یادگیری ماشین به اندازه ی کافی پیچیده نباشد که بتواند روابط معنا دار میان ویژگی ها و متغیر هدف را به خوبی یاد بگیرد.
- در این حالت مقدار Bias بالا و واریانس پایینی داریم.
- مدلی که دارای Underfitting باشد، به خوبی آموزش نمی بیند و هم بر روی داده های train و هم test به خوبی عمل نمی کند.



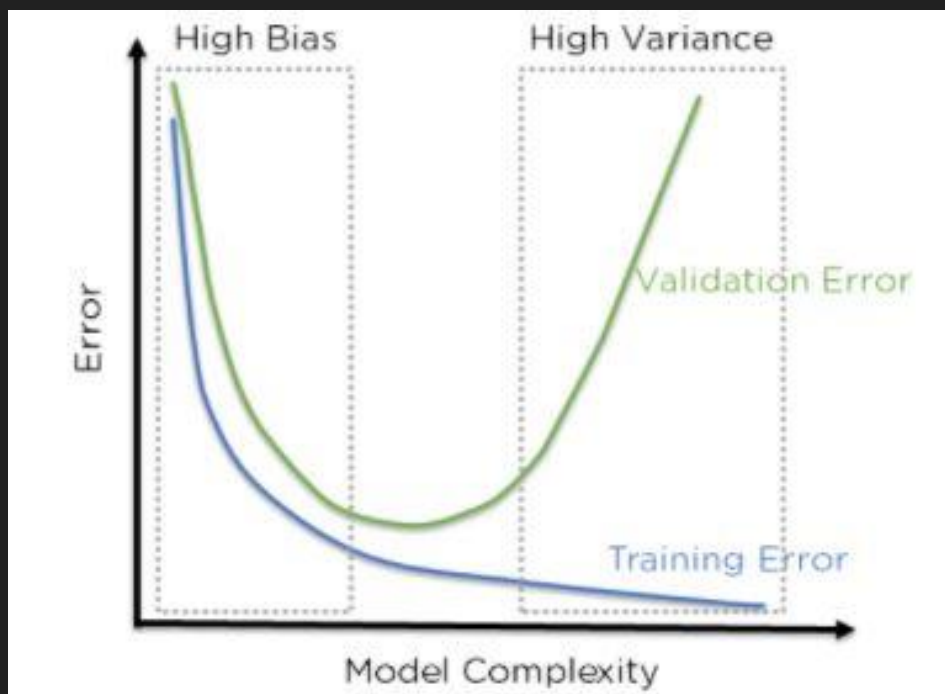


# روش های مقابله با Underfitting

- کاهش تنظیم (Decreasing Regularization): با کاهش میزان Regularization پیچیدگی و تنوع مدل بیشتر می شود و امکان آموزش بهتر فراهم می شود.
- افزایش زمان آموزش بر روی داده
- انتخاب ویژگی ها به طور موثر



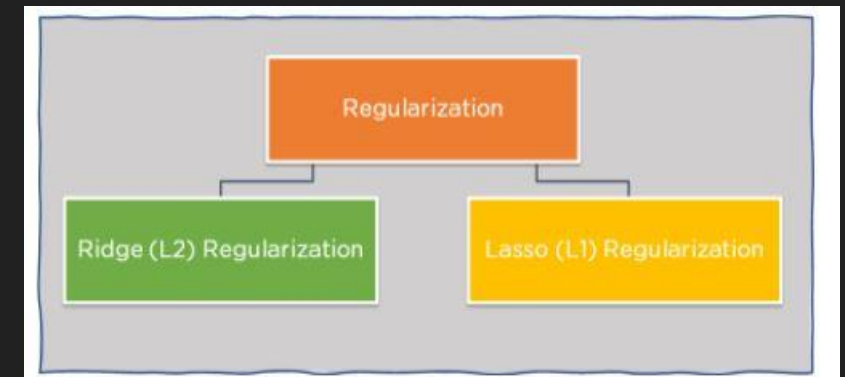
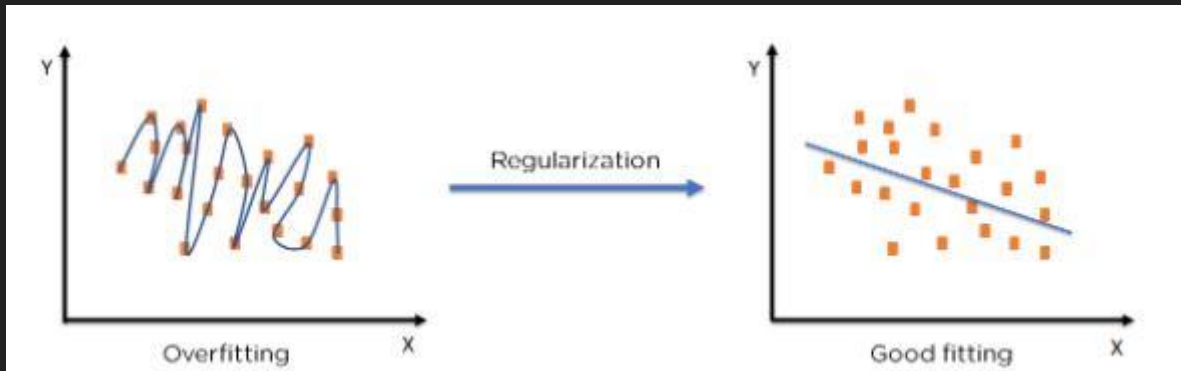
# بایاس و واریانس



- بایاس زمانی اتفاق می افتد که مدل یادگیری ماشین، قابلیت انعطاف پذیری کمی برای یادگیری داشته باشد. چنین مدلی، توجه زیادی به داده ها می کند و قوانین را بیش از حد ساده سازی می کند. در نتیجه، میزان خطا بر روی داده های train و prediction افزایش می یابد.
- واریانس میزان حساسیت مدل یادگیری ماشین بر روی داده ها را مشخص می کند. یک مدل یادگیری ماشین با واریانس بالا، توجه زیادی به داده های آموزشی می کند و نمی تواند به خوبی قوانین را کلی سازی یا Generalize کند در نتیجه خطای یادگیری بر روی prediction بالا است. (خطای یادگیری در این حالت بر روی train کم است. چون داده های train را خیلی خوب یاد می گیرد یا بهتر بگوییم، بیش از اندازه خوب یاد می گیرد).
- پس میزان بایاس و واریانس باید به درستی انتخاب شوند که در کنار یکدیگر به فرآیند یادگیری کمک کنند.

# نظم دهی یا Regularization

- این روش باعث اعمال محدودیتی هایی بر روی مدل یادگیری ماشین می شود تا از Overfitting یا Underfitting مدل جلوگیری کند.
- در نظم دهی نوع  $L_1$  یا  $L_2$  می توانیم برای تابع زیان یا همان Loss function یک مجازات در نظر بگیریم تا ضرایب Coefficient ها را به سمت صفر سوق دهد. نظم دهی نوع  $L_2$  اجازه می دهد وزن ها به سمت صفر پیش بروند اما صفر نشوند. در حالیکه در  $L_1$  وزن ها به صفر می رسد.



# L1 Regularization یا Lasso Regression یا نُرم ۱

- تصور کنید الگوریتمی نیاز داریم تا رتبه کنکور فردی را بر اساس سوابق او پیشبینی کند. طبیعی است که همه ویژگی‌های در سابقه یک فرد تاثیر یکسانی در کسب رتبه ندارند. مثلاً معدل نمرات فرد در پیشبینی بیشتر به کار می‌آید تا سوابق فعالیت‌های داوطلبانه و خارج از درس، و یا میانگین رتبه‌ها در منطقه آن فرد هم موثرتر از شاخص فیزیکی BMI اوست. پس با کمک نُرم L1 در حین آموزش به ویژگی‌های کمتر موثر وزن بسیار بسیار کوچکی (نزدیک به صفر) اعمال می‌شود چرا که تاثیر آن‌ها کمتر است. به طور مثال وزن مربوط به میانگین رتبه‌ها در منطقه به سمت عددی غیر صفر میل می‌کند، اما وزن شاخص BMI فرد مدام کوچکتر می‌شود تا به صفر برسد.
- قسمت قرمز رنگ، مقدار پینالتی یا جریمه را مشخص می‌کند. در این قسمت، مجموع قدر مطلق همه‌ی وزن‌ها در یک مقدار ثابت (لاندا) ضرب می‌شود. با تغییر لاندا، الگوریتم سعی دارد تا میزان تابع هزینه را کمینه کند.

Cost function:

$$L(x, y) \equiv \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n |\theta_i|$$

# L۲ Regularization یا Ridge Regression یا نُرم ۲

- در Ridge، تابع هزینه با توجه به مجموع مجذور وزن ها بدست می آید.
- در Ridge، مانند Lasso، سعی الگوریتم بر کوچک کردن وزن ها هست، اما برخلاف Lasso، در این حالت وزن ها را صفر نمی کند.
- هنگامی که داده های Outlier در دیتاست موجود باشد، این روش به خوبی عمل نمی کند چرا که در نقاط Outlier خطای پیش بینی مدل بسیار زیاد است و با داشتن جمله ی پینالیتی مانند زیر، وزن ها کوچکتر خواهد شد.
- مدل یادگیری ماشینی که از Ridge استفاده می کند، زمانی بهتر عمل می کند که تمامی ویژگی های داده ورودی روی پیش بینی هدف یا خروجی تاثیر گذار بوده و همچنین، وزن های داخل مدل به طور مساوی مقداردهی اولیه شده باشند.

Cost function:

$$L(x, y) \equiv \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n \theta_i^2$$

# مقایسه L۱ و L۲

## ○ نُرم L۱

- مجموع قدر مطلق وزن ها را به هزینه/خطا اضافه می کند.
- نُرم L۱ یک مدل خلوت تر ایجاد می کند.
- ممکن است براساس اثردهی زیرمجموعه های متفاوت از ویژگی ها، مدل با نُرم L۱ چند یادگیری متفاوت ایجاد کند.
- با نُرم L۱ خاصیت feature selection حاصل می شود.
- نُرم L۱ در مواجهه با نقاط outlier در داده های ورودی، عملکرد بهتری و مقاومتری دارد.
- نُرم L۱ می تواند مدلی ساده و قابل تفسیر ایجاد کند اما همین سبب می شود تا الگوریتم قادر نباشد الگوهای پیچیده را یاد بگیرد.

## ○ نُرم L۲

- مجموع مجذور وزن ها را به هزینه/خطا اضافه می کند.
- نُرم L۲ یک مدل خلوت ایجاد نمی کند.
- نُرم L۲ تنها یک یادگیری/راه حل در الگوریتم یادگیری ماشین ایجاد می کند و براساس استفاده از زیرمجموعه های متفاوت از ویژگی ها عمل نمی کند.
- با نُرم L۲ خاصیت feature selection حاصل نمی شود.
- نُرم L۲ در مواجهه با نقاط outlier در داده های ورودی، عملکرد خوبی ندارد.
- زمانی که مقدار هدف پیش بینی تابعی از همه ویژگی ها داده ورودی باشد، نُرم L۲ یادگیری بهتری را سبب می شود.
- با نُرم L۲ بر خلاف نُرم L۱، می توان الگوهای پیچیده را در داده ورودی یاد گرفت.

# ترکیب دو روش قبل

○ آیا بهتر نیست هر دو را با همدیگر استفاده کنیم تا ضعف های یکدیگر را پوشش دهند؟ آیا منجر به نتیجه بهتر خواهد شد؟ جواب مثبت است. به الگوریتم Elastic net regularization که ترکیبی از هر دو مورد  $L_1$  و  $L_2$  regularization را استفاده می کند نگاهی بیاندازید!

# پایان

با تشکر از توجه تان، اوقات خوشی را برایتان آرزومندم.