

One hot encoding زمانی استفاده میشود که داده های دسته ای یا کتوگوریکال را به داده باینری تبدیل کنیم، از معایب آن زمانی که دسته ها زیاد باشد باعث افزایش فیچر های مورد نظر و پیچیدگی مدل میشود

با استفاده از تبدیل log and exponential transformation میتوان داده های با توزیع log normal را به normal تبدیل کرد.

روش bootstrapping از روشهای نمونه گیری مجدد است و روش کارآمدی برای محاسبه میزان دقت و خطای استاندارد (standard error) متغیر تخمین زده شده است. در روش bootstrapping به صورت تصادفی  $n$  عضو با جایگذاری انتخاب میکنیم برای مثال فرض کنید مقدار معینی پول داریم که  $m$  یخواهیم بخش  $\alpha$  درصد آن را در یک پروژه و  $1-\alpha$  درصد آن را در یک پروژه دیگر سرمای گذاری کنیم. اگر میزان سود پروژه اول  $X$  و میزان سود پروژه دوم  $Y$  باشد، به دنبال یافتن  $\alpha$  ای هستیم که به ازای آن، میزان ریسکمان یعنی  $Var(\alpha X + (1-\alpha)Y)$  کمینه شود

در روش cross validation هر بار یک زیرمجموعه از داده های آموزشی را کنار م یگذاریم، مدل را با استفاده از بقیه داده ها آموزش میدهم و از داده های کنار گذاشته شده به عنوان داده های تست استفاده میکنیم.

در 2-fold cross-validation داده ها را بصورت رندوم به 2 بخش هم اندازه تقسیم می کنیم. یک بخش به عنوان داده های ارزیابی و دیگری بعنوان داده های آموزشی استفاده می شود و بار دیگر همین کار را با بخش دیگر انجام می دهیم در نتیجه دو تا معیار خطای تست خواهیم داشت که می توانیم از آن ها میانگین بگیریم و بعنوان خطای تست مدل گزارش دهیم. مزیت: هر بخش یکبار بعنوان داده های ارزیابی استفاده شده است و از همه ی مشاهدات استفاده کرده ایم.

در 5\*2 fold cross-validation، 5 بار 2-fold cross-validation را اجرا می کنیم. می توان خطای تست در این مورد با میانگین گرفتن از 5 نتیج هی 2 fold cross-validation بدست آورد. این روش که در واقع تکرار روش k-fold cv است باعث می شود تا تخمین بهتری از خطای تست بدست آید

براساس means-k و مفهوم متغیر inertia موجود در آن که برابر مجموع مربعات فاصله ها از نزدیک ترین مرکز خوشه است. وقتی نمودار مقادیر  $k$  بر حسب inertia را رسم شود از مقدار  $k$  مشخص به

بعد نمودار به صورت خطی کم می شود. ولی شرایط قبل نمایی است. به این نقطه elbow گفته میشود. و مقدار بهینه برای  $k$  این نقطه است. در نمودار زیر نیز چنین الگویی وجود دارد. در نمودار از مرتبه ۳ تا ۵ مقدار خطی است و بنابراین بهترین مدل برای پیچیدگی پیدا میشود. رفتار تغییر بایاس همیشه به صورت خطی نیست اما با مشاهده واریانس و افزایش واریانس نسبت به کاهش بایاس می توان نقطه بهینه برای انتخاب پیچیدگی مدل را پیدا کرد پس نظر به طور کلی نمیتوان گفت.