# Index

## Z
zero-order optimization, 21